

# Video-MMMU: Evaluating Knowledge Acquisition from Multi-Discipline Professional Videos

Kairui Hu<sup>1</sup> Penghao Wu<sup>1</sup> Fanyi Pu<sup>1</sup> Wang Xiao<sup>1</sup> Yuanhan Zhang<sup>1</sup> Xiang Yue<sup>2</sup>  
 Bo Li<sup>1</sup> Ziwei Liu<sup>1\*</sup>

<sup>1</sup>S-Lab, Nanyang Technological University

<sup>2</sup>Carnegie Mellon University

<https://videommmu.github.io/>

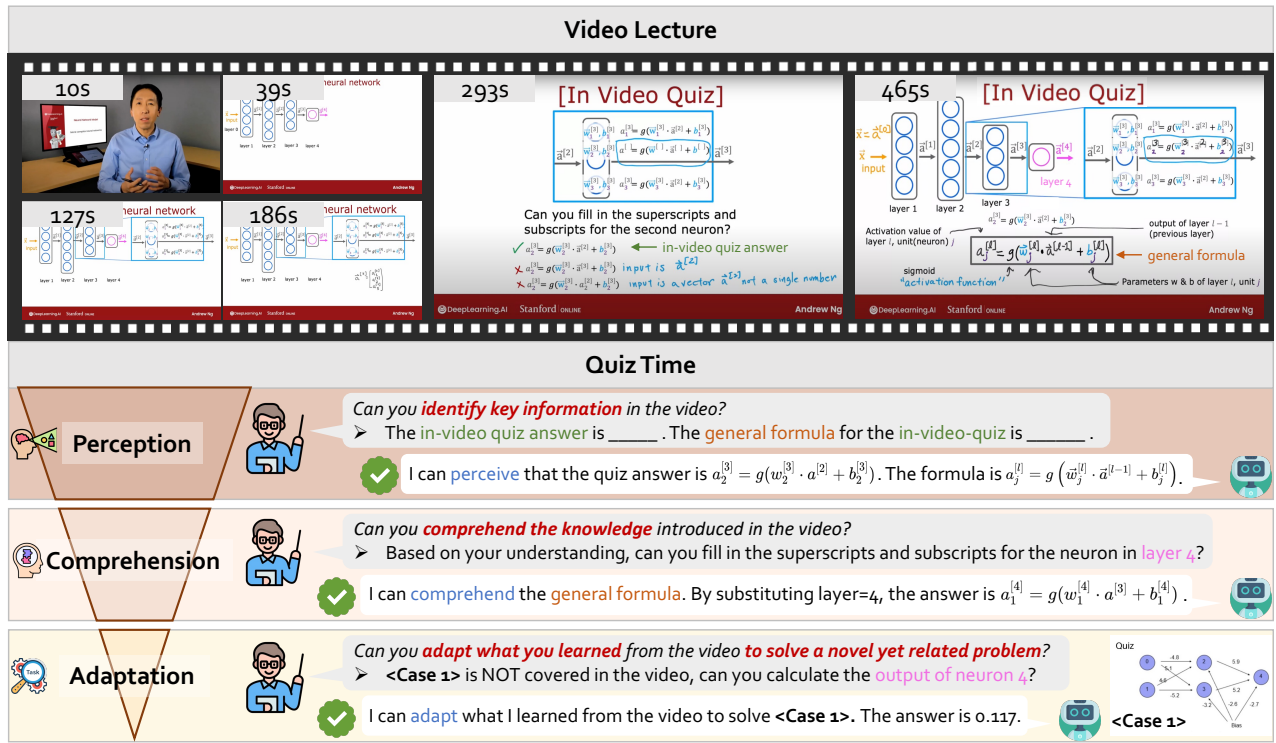


Figure 1. An illustration of **Video-MMMU**: Evaluating the knowledge acquisition capability from videos through three cognitive stages: **1) Perception**: if models can identify key information related to knowledge; **2) Comprehension**: if models can interpret the underlying concepts; **3) Adaptation**: if models can adapt the knowledge from videos to novel scenarios.

## Abstract

Humans acquire knowledge through three cognitive stages: perceiving information, comprehending knowledge, and adapting knowledge to solve novel problems. Videos serve as an effective medium for this learning process, facilitating a progression through these cognitive stages. However, existing video benchmarks fail to systematically evaluate the knowledge acquisition capabilities in Large Multimodal

Models (LMMs). To address this gap, we introduce Video-MMMU, a multi-modal, multi-disciplinary benchmark designed to assess LMMs’ ability to acquire and utilize knowledge from videos. Video-MMMU features a curated collection of 300 expert-level videos and 900 human-annotated questions across six disciplines, evaluating knowledge acquisition through stage-aligned question-answer pairs: Perception, Comprehension, and Adaptation. A proposed knowledge gain metric,  $\Delta_{\text{knowledge}}$ , quantifies improvement in performance after video viewing. Evaluation of LMMs reveals a

*steep decline in performance as cognitive demands increase and highlights a significant gap between human and model knowledge acquisition, underscoring the need for methods to enhance LMMs’ capability to learn and adapt from videos.*

## 1. Introduction

Humans acquire knowledge through three fundamental cognitive stages outlined in Bloom’s taxonomy [8]: **1) perceiving** information, **2) comprehending** knowledge, and **3) adapting** knowledge to solve novel problems. Video serves as an ideal medium for this learning process, enabling a natural progression from information intake to practical application, making video-based learning a valuable tool for knowledge acquisition [11, 30, 43]. Consider learning neural network forward propagation through video lectures (Fig. 1): learners first recognize fundamental concepts like activation functions, then demonstrate understanding through exercises, and ultimately apply this knowledge to solve novel exam problems. This progression naturally aligns with Bloom’s cognitive stages, providing a systematic framework for assessing knowledge acquisition from videos.

For Large Multimodal Models (LMMs) to operate effectively in the wild like humans, learning from videos is an essential capability for continuous knowledge acquisition. However, existing video benchmarks lack systematic evaluation of this critical ability. To bridge this gap, we introduce **Video-MMMU**, a massive multi-modal, multi-disciplinary video benchmark that evaluates the knowledge acquisition capability from educational videos through three main features: **1) Knowledge-intensive Video Collection:** Our dataset comprises 300 expert-level videos spanning 6 professional disciplines: Art, Business, Science, Medicine, Humanities, and Engineering, with 30 subjects distributed among them. **2) Knowledge Acquisition-based Question Design:** Each video includes three question-answer pairs aligned with the three knowledge acquisition stages: Perception (identifying key information related to the knowledge), Comprehension (understanding the underlying concepts), and Adaptation (applying knowledge to new scenarios). **3) Quantitative Knowledge Acquisition Assessment:** We propose a knowledge acquisition metric, denoted as  $\Delta_{\text{knowledge}}$ , to measure performance gains on practice exam questions after learning from videos. This metric enables us to quantitatively evaluate how effectively large multimodal models (LMMs) can assimilate and utilize the information presented in the videos to solve real-world, novel problems.

We evaluate both open-source and proprietary LMMs on Video-MMMU, revealing several key findings: **1) Progressive Performance Decline:** Model performance decreases as cognitive demands increase. While models perform relatively better on perception tasks, their accuracy drops notably on comprehension tasks and declines further on adap-

tation tasks. **2) Knowledge Acquisition from videos is Challenging:** The knowledge acquisition metric  $\Delta_{\text{knowledge}}$  reveals a significant gap between human and model performance. While humans achieve substantial improvement ( $\Delta_{\text{knowledge}} = 33.1\%$ ) after watching the videos, even the top performing models show smaller knowledge gains (GPT-4o [27]:  $\Delta_{\text{knowledge}} = 15.6\%$ , Claude-3.5-Sonnet [1]:  $\Delta_{\text{knowledge}} = 11.4\%$ ). This limitation underscores a challenge in current LMMs. While humans naturally acquire knowledge through video-based learning, having developed this capability through classroom learning and educational experiences throughout life, LMMs struggle to effectively learn from videos. These findings emphasize the need for further research to enhance how LMMs acquire and utilize video-based information, bringing them closer to human-level learning processes.

## 2. Related Work

### 2.1. VideoQA Benchmarks

Existing video benchmarks focus primarily on visual understanding tasks, including action understanding [14, 22, 28, 38, 39, 44], temporal reasoning [3, 18, 20, 31, 34, 37, 42], and video captioning [4, 35, 40, 41, 53]. Several benchmarks enhance scene interpretation by incorporating external knowledge, including KnowIT-VQA [10] and WorldQA [50]. Recent benchmarks like Video-MME [9], MMBench-Video [7], and MLVU [52] have expanded the scope to assess multi-tasking and multi-domain video understanding. While these benchmarks recognize videos as visual scenes for interpretation, Video-MMMU uniquely recognizes video as an educational medium, emphasizing knowledge-driven question-answering on videos.

### 2.2. Knowledge-driven Benchmarks

As AI systems progress toward Expert AGI [24], knowledge-driven benchmarks have emerged to evaluate models’ professional expertise. Early benchmarks such as AGIEval [51] and ARC [2] focus on standardized exams and science questions, respectively. MMLU [13] expands evaluation across STEM disciplines, while MMLU-Pro [36] introduces more challenging reasoning-focused questions. Multi-modal benchmarks extend this evaluation scope further. ScienceQA [21] assesses multi-modal reasoning on elementary to high-school science questions. MMMU [45] advances to college-level questions requiring subject-specific knowledge and deliberate reasoning. MMMU-Pro [46] enhances MMMU questions for more robust evaluation. While these benchmarks evaluate models’ pre-trained knowledge and reasoning abilities on text and images, Video-MMMU uniquely focuses on assessing how effectively models can acquire and apply knowledge from videos.


Art	Humanities	Medicine
 <p><b>Question:</b> What does the speaker say when introducing Peter Paul Rubens at the end of the video? Select the option that precisely matches the speaker's statement.  <b>Options:</b>            (A) Peter Paul Rubens was a famous Baroque...            (B) Peter Paul Rubens is regarded as a prolific artist...            .....  <b>(I) Peter Paul Rubens was the most important...</b>            (J) Peter Paul Rubens is celebrated for his dynamic...</p> <p><b>Track:</b> Perception, <b>Video Type:</b> Concept-introduction video, <b>Subject:</b> Art Theory, <b>QA Type:</b> Automatic Speech Recognition (ASR)</p>	 <p><b>Question:</b> Based on your understanding of cultural universals from the video, determine which of the following statements are correct:            Statement 1: All human cultures have some...            Statement 2: The video uses the example of...            Statement 3: At 3:35, the video implies that ...            Statement 4: ... Statement 5: ...  <b>Options:</b>            (A) Statement 1 (B) Statement 2,3 <b>(C) Statement 3,4</b>            (D) Statement 2,4,5 .....(J) Statement 2,4</p> <p><b>Track:</b> Comprehension, <b>Video Type:</b> Concept-introduction video, <b>Subject:</b> Sociology, <b>QA Type:</b> Concept Comprehension (CC)</p>	 <p><b>Question:</b> Can you identify the abnormality on this plain film of the pelvis? &lt;image 1&gt;  <b>Options:</b>            (A) Bone cyst            (B) Acute hip fracture            (C) Osteoarthritis            (D) Surgical hardware  <b>(E) Resection of the pubic symphysis</b>            ...            (J) Bone infection</p>  <p><b>Track:</b> Adaptation, <b>Video Type:</b> Concept-introduction video, <b>Subject:</b> Clinical Medicine, <b>QA Type:</b> Case Study Analysis (CSA)</p>
Business	Science	Engineering
 <p><b>Question:</b> According to the video, a minimum price control on alcoholic drinks is intended to reduce consumption from <math>Q_1</math> to _____, addressing negative externalities. The policy raises the price to _____ above the free market price of _____. Fill in the blanks based on the video content.  <b>Options:</b>  <b>(A) <math>Q^*</math>, <math>P_{min}</math>, <math>P_1</math></b> (B) <math>Q^*</math>, <math>P_1</math>, <math>P_{min}</math> (C) <math>Q_1</math>, <math>P_{min}</math>, <math>P_2</math>            (D) <math>Q_2</math>, <math>P_1</math>, <math>P_{min}</math> (E) <math>Q^*</math>, <math>P_2</math>, <math>P_1</math> ... (F) <math>Q_1</math>, <math>P_2</math>, <math>P_{min}</math>            (G) <math>Q_2</math>, <math>P_{min}</math>, <math>P_1</math>. (H).... (I).... (J) <math>Q_1</math>, <math>P_1</math>, <math>P_{min}</math></p> <p><b>Track:</b> Perception, <b>Video Type:</b> Problem-solving video, <b>Subject:</b> Economics, <b>QA Type:</b> Optical Character Recognition (OCR)</p>	 <p><b>Question:</b> In the video, Example Question (1) is solved with an angle <math>\theta=25</math> degrees. If the angle <math>\theta</math> is adjusted to 30 degrees while all other conditions remain unchanged, what will be the updated result for Example Question (1) as explained in the video?  <b>Options:</b>            (A) 4.00 seconds (B) 2.82 seconds (C) 3.50 seconds            (D) 2.50 seconds <b>(E) 3.04 seconds</b> (F) 2.00 seconds            (G) 3.15 seconds (H) 1.85 seconds (I) 2.25 seconds            (J) 3.85 seconds</p> <p><b>Track:</b> Comprehension, <b>Video Type:</b> Problem-solving video, <b>Subject:</b> Math, <b>QA Type:</b> Problem-solving Strategy Comprehension (PSC)</p>	 <p><b>Question:</b> Based on what you learned from the video, write the Fourier series for the three voltage waveforms in (a) of &lt;image 1&gt;.  <b>Options:</b>            (A) <math>(4/\pi)(\sin(\pi t)+(1/2)\sin(3\pi t)+(1/4)\sin(5\pi t)+\dots)</math>  <b>(B) <math>(4/\pi)(\sin(\pi t)+(1/3)\sin(3\pi t)+(1/5)\sin(5\pi t)+\dots)</math></b>            (C) <math>(4/\pi)(\sin(\pi t)+(1/2)\sin(2\pi t)+(1/4)\sin(4\pi t)+\dots)</math>            .....            (J) <math>(4/\pi)(\sin(\pi t)+(1/4)\sin(3\pi t)+(1/6)\sin(5\pi t)+\dots)</math></p> <p><b>Track:</b> Adaptation, <b>Video Type:</b> Problem-solving video, <b>Subject:</b> Electronics, <b>QA Type:</b> Problem-solving Strategy Adaptation (PSA)</p>

Figure 2. Sampled Video-MMMU examples across 6 academic disciplines and 3 tracks. The examples are organized in two rows based on distinct video types: (1) Concept-Introduction videos (top row) focus on teaching factual knowledge, fundamental concepts, and theories through explanatory content, while (2) Problem-Solving videos (bottom row) demonstrate step-by-step solutions to an example question.

### 3. Video-MMMU Dataset

We introduce Video-MMMU (Massive Multi-discipline Multimodal Understanding), a video benchmark designed to evaluate knowledge acquisition from educational videos across 30 subjects in 6 professional disciplines: Art, Business, Medicine, Science, Humanities, and Engineering. The video distribution across disciplines is shown in Fig. 3a.

#### 3.1. Video Collection

The dataset consists of 300 college-level educational videos, systematically curated through a rigorous three-phase process: **1) Topic Selection:** Domain experts conduct a comprehensive analysis of college curricula across 30 subjects, establishing a diverse pool of 450 foundational assessment topics. **2) Video Curation:** Leveraging GPT-4o [27], we generated 10 search queries per topic. These search queries are processed through the YouTube Data API to create an

initial candidate video pool. **3) Quality Assurance:** We implemented a three-tier review protocol: First, annotators cross-check to filter out videos with poor audio-visual quality or irrelevant content. Second, we employ GPT-4o [27] to assess the technical depth of the videos by analyzing 10 sampled frames from each video. We prioritize in-depth lectures, tutorials, and detailed problem-solving demonstrations while excluding beginner-level introductions and superficial overviews. Finally, domain experts verify alignment with college curriculum standards and confirm appropriate domain knowledge depth.

The Video-MMMU dataset comprises two distinct categories: **1) Concept-introduction Videos:** These videos provide comprehensive explanations of factual knowledge, including fundamental concepts and theories. **2) Problem-solving Videos:** These videos demonstrate step-by-step problem solutions, particularly in STEM disciplines where

systematic reasoning and detailed calculations are required.

## 3.2. QA Annotation

### 3.2.1. QA Taxonomy

We annotate questions across three cognitive stages: Perception, Comprehension, and Adaptation, each assessing progressively deeper levels of knowledge acquisition.

**Perception Questions** assess the ability to perceive information from videos through: **1) Optical Character Recognition (OCR):** These questions require identifying and extracting key details from visual content, including formulas, data points, charts, and handwritten notes. An example is shown in Fig. 2 (Business), where the question requires extracting multiple economic variables from handwritten notes. **2) Automatic Speech Recognition (ASR):** These questions assess the ability to accurately transcribe spoken content into text, as illustrated in Fig. 2 (Art).

**Comprehension Questions** evaluate the ability to understand knowledge presented in videos through: **1) Concept Comprehension (CC):** These questions assess understanding of concepts introduced in the videos. We primarily use a multiple-answer multiple-choice (MAMC) format, where each question presents 4-10 statements about video content, with multiple correct statements possible. As shown in Fig. 2 (Humanities), one must identify all correct statements about the video content to demonstrate a comprehensive understanding. **2) Problem-solving Strategy Comprehension (PSC):** For videos demonstrating step-by-step solutions to example questions, an intuitive way to assess the understanding of the solution is to test the same question with different input values. As illustrated in Fig. 2 (Science), when a video demonstrates trajectory time calculation with a 25-degree angle, the question changes this to 30 degrees. This approach verifies comprehension of the underlying solution strategy rather than the memorization of answers. The cognitive difficulty lies between perception and adaptation, requiring new calculations while following the same reasoning process in the video.

**Adaptation Questions** assess the ability to adapt video knowledge to new scenarios: **1) Case Study Analysis (CSA):** These questions evaluate the application of concepts to novel real-world scenarios. As shown in Fig. 2 (Medicine), while the video explains various pelvic pathologies, the question requires analysis of a new patient’s pelvic radiograph to identify specific abnormalities. This tests the model’s ability to adapt theoretical knowledge from videos to practical clinical diagnosis. **2) Problem-solving Strategy Adaptation (PSA):** These questions evaluate how learners adapt learned solution methods to new problems. For instance, in Fig. 2 (Engineering), the video demonstrates the calculation of Fourier series for one type of waveform, while the question presents a different waveform pattern. To solve this new problem, one needs to identify key similarities and

differences between the video example and the new problem, then adjust the solution method accordingly. The distribution of these question types is illustrated in Fig. 3b.

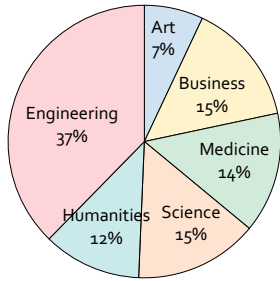
### 3.2.2. Annotation and Quality Control

**Annotation Process:** Our annotation follows a multi-stage process to ensure quality: **1) Initial Annotation:** Annotators thoroughly review each video and annotate three questions aligned with our cognitive tracks, following the QA taxonomy shown in Fig. 3b. To enhance assessment rigor, we annotate 10 options for each multiple-choice question (MCQ). **2) Quality Assurance:** Firstly, annotators cross-check each other’s questions for consistency and clarity. Secondly, QA pairs are processed by OpenAI o1 [26] to refine the language and verify the correctness of ground-truth answers. Thirdly, domain experts review each question for technical accuracy and alignment with the intended cognitive stages. For Adaptation questions, experts verify that the question tests the same knowledge presented in the video but in a novel scenario, ensuring they utilize the same concepts, formulas, or similar problem-solving strategies. Finally, we employ Gemini 1.5 Pro [32] to analyze each video-question pair and determine whether audio might be helpful to solve the question, as shown in Fig. 3c. This analysis will benefit more future Large Multimodal Models (LMMs) with audio processing capabilities.

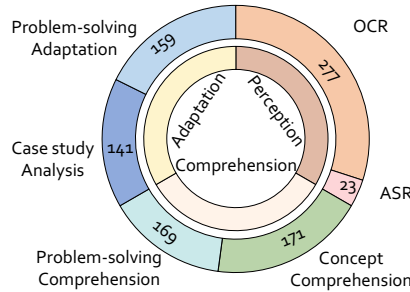
**Question Sources:** For the Perception and Comprehension tracks, questions are manually created by our annotators. For the Adaptation track, which requires practical problems from exams and case studies, our approach varies by discipline. In Science, Engineering, Medicine, and Business, we source questions from MMMU [45] and MMMU-pro [46], which provide validated college exam questions well suited for testing knowledge adaptation. For Art and Humanities, where adaptation requires more context-dependent assessment, we manually create case study questions to ensure alignment with video concepts.

## 3.3. Comparison with Existing Benchmarks

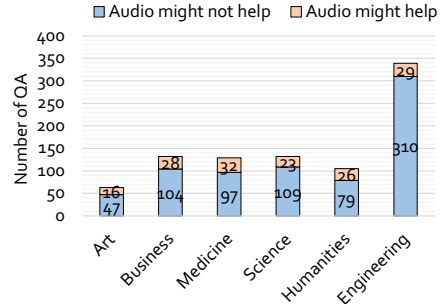
Video-MMMU distinguishes itself through its emphasis on how models can learn and apply knowledge from professional educational videos. Our videos feature comprehensive lectures, tutorials, and step-by-step problem-solving demonstrations, delivering dense information through multiple visual formats, including charts, diagrams, and handwritten explanations. With an average duration of 506.2 seconds, the videos provide extensive coverage of domain-specific knowledge across various disciplines. As shown in Table 1, our questions are substantially longer than existing benchmarks, averaging 75.7 words per question, reflecting the complexity of knowledge-driven evaluation. We systematically evaluate knowledge acquisition from videos through three cognitive stages. The Adaptation track advances video-based learning evaluation beyond basic content understanding to assess



(a) Video distribution across disciplines.



(b) QA distribution across types.



(c) QA distribution with respect to audio.

Figure 3. Taxonomy of QA types and video disciplines.

Benchmarks	Video Domain	Question Length	Video Duration	Knowledge driven
Video-MME [9]	Open	35.7	1017.9	✗
MMBench-Video [7]	Open	10.9	165.4	✗
Video-Bench [25]	Open	21.3	56.0	✗
TempCompass [20]	Open	49.2	11.4	✗
MVBench [17]	Open	27.3	16.0	✗
AutoEval-Video [5]	Open	11.9	14.6	✗
Video-MMMU	Professional	75.7	506.2	✓

Table 1. Comparison of Video-MMMU and other widely adopted video benchmarks.

how effectively models can apply the acquired knowledge to novel problems.

## 4. Experiments

### 4.1. Settings

**Baselines.** We evaluate open-source LMMs including LLaVA-OneVision [15], LLaVA-Video [49], LongVA [48], VILA-1.5 [19], Qwen2-VL [33], InternVL2 [6], Llama-3.2 [23], MAmmoTH-VL [12], Aria [16]; and proprietary models GPT-4o [27], Gemini 1.5 Pro [32], Gemini 1.5 Flash [32], Claude-3.5-Sonnet [1]. The numbers of sampled frames are 32 for LLaVA-OneVision, 64 for LLaVA-Video, 64 for LongVA, 32 for VILA-1.5, 32 for InternVL2, 10 for Llama-3.2, 32 for MAmmoTH-VL, 64 for Aria, 20 for Claude-3.5-Sonnet, 50 for GPT-4o.

**Human Experts.** To assess the performance of Human Experts, we recruited senior undergraduate students and instructed them to complete the following tests: The students first attempted the Adaptation question without viewing the videos. Subsequently, they watched each assigned video and answered the corresponding Perception, Comprehension, and Adaptation questions. While students could refer to course materials and notes, they were not allowed to search for answers on the Internet.

**Inputs.** We provide videos and questions as inputs for the Perception and Comprehension tracks. For the Adaptation track, we append the question’s image to the end of each video. We add a prompt to indicate that the image for the Adaptation track question appears in the final frame.

**Evaluations.** We evaluate model outputs using an automated, rule-based pipeline. The system employs regular expressions to extract key elements such as option letters and numerical values. Responses lacking valid answers are marked as incorrect. We use the micro-averaged accuracy as our evaluation metric. The evaluation is conducted using LMMs-Eval [47].

## 4.2. Main Results

### 4.2.1. Performance by Track

**Human vs. Model Performance:** Human experts outperform models across all tracks, with Claude achieving the highest model scores but still showing a gap to humans. Both humans and models exhibit declining performance from perception through comprehension to adaptation, indicating that deeper cognitive stages require more advanced capabilities. **Perception Track:** Many models achieve an accuracy over 50%, suggesting perception is a more fundamental capability among the three stages. **Comprehension Track:** Comprehending college-level knowledge from videos requires pre-trained knowledge as a foundation. Compared to the Perception score, most open-source models show a 10 ~ 20% decline in Comprehension score, while proprietary models show less performance decline and generally achieve higher comprehension scores, demonstrating their superior capabilities in comprehending knowledge-intensive videos. **Adaptation Track:** Adaptation emerges as the most challenging stage, with most models scoring below 50%. Even top-performing models like Claude-3.5-Sonnet exhibit a substantial performance decline in Adaptation. This indicates a natural gap between theoretical understanding and practical application. While models might understand the knowledge from videos at a surface level, they currently lack

Model	Overall	Results by Track			Results by Discipline					
		Perception	Comprehension	Adaptation	Art.	Biz.	Sci.	Med.	Hum.	Eng.
Random Choice	14.00	12.00	14.00	16.00	11.11	12.88	12.12	22.48	10.48	13.57
Human Expert	74.44	84.33	78.67	60.33	80.95	78.79	74.24	70.54	84.76	69.91
<b>Proprietary LMMs</b>										
Gemini 1.5 Flash [32]	49.78	57.33	49.00	43.00	63.49	53.03	43.18	49.61	59.05	45.72
Gemini 1.5 Pro [32]	53.89	59.00	53.33	49.33	57.14	59.09	49.10	57.42	58.10	50.31
GPT-4o [27]	61.22	66.00	62.00	<b>55.67</b>	69.52	66.88	51.55	<b>64.76</b>	69.52	57.13
Claude-3.5-Sonnet [1]	<b>65.78</b>	<b>72.00</b>	<b>69.67</b>	<b>55.67</b>	66.67	<b>75.00</b>	<b>56.06</b>	58.14	<b>75.24</b>	<b>66.08</b>
<b>Open-source LMMs</b>										
VILA1.5-8B [19]	20.89	20.33	17.33	25.00	34.92	14.39	19.70	19.38	21.91	21.53
LongVA-7B [48]	23.98	24.00	24.33	23.67	41.27	20.46	21.97	24.03	23.81	23.01
Llama-3.2-11B [23]	30.00	35.67	32.33	22.00	39.68	28.79	21.21	35.66	33.33	28.91
LLaVA-OneVision-7B [15]	33.89	40.00	31.00	30.67	49.21	29.55	34.85	31.78	46.67	29.20
VILA1.5-40B [19]	34.00	38.67	30.67	32.67	57.14	27.27	23.49	37.99	41.91	32.45
LLaVA-Video-7B [49]	36.11	41.67	33.33	33.33	65.08	34.09	32.58	42.64	45.71	27.43
InternVL2-8B [6]	37.44	47.33	33.33	31.67	55.56	34.09	30.30	34.11	41.91	38.05
MAmmoTH-VL-8B [12]	41.78	51.67	40.00	33.67	47.62	37.88	36.36	36.43	49.52	43.95
LLaVA-OneVision-72B [15]	48.33	59.67	42.33	43.00	61.91	46.21	40.15	54.26	60.00	43.95
LLaVA-Video-72B [49]	49.67	59.67	46.00	43.33	69.84	44.70	41.67	58.92	57.14	45.13
Aria [16]	50.78	65.67	46.67	40.00	<b>71.43</b>	47.73	44.70	58.92	62.86	43.66

Table 2. Video-MMMU Evaluation Results across three cognitive tracks (Perception, Comprehension, Adaptation) and six disciplines (Art, Business, Science, Medicine, Humanities, Engineering).

the advanced capability to effectively acquire and apply what they learned from the video to solve practical problems.

#### 4.2.2. Performance by Discipline

Model performance varies across disciplines. Models demonstrate superior performance in Art and Humanities disciplines, where videos primarily focus on conceptual presentation. In comparison, they achieve lower accuracies in Science, Engineering, Business, and Medicine, which demand quantitative reasoning and interpretation of detailed technical visuals such as diagrams and handwritten notes. This performance differential suggests models are generally more adept at processing factual knowledge but underperform in domains requiring complex computation, deliberate reasoning, and visual analysis.

#### 4.3. Impact of Audio Transcript

Audio conveys information in knowledge-intensive videos. To study the impact of audio transcripts, we use OpenAI Whisper [29] to generate audio transcripts and append them to the input prompt. We conduct evaluation on the top-performing open-source model Aria [16] and proprietary model Claude-3.5-Sonnet [1].

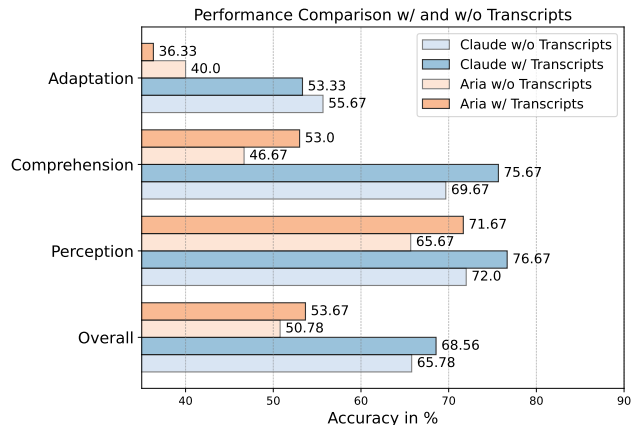
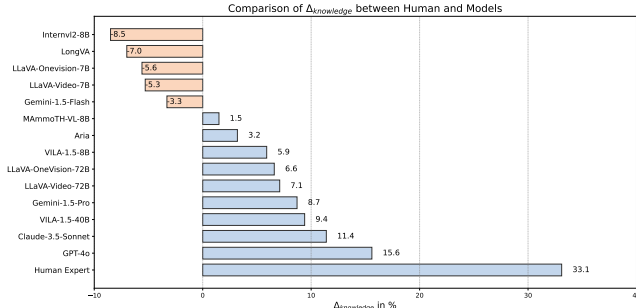
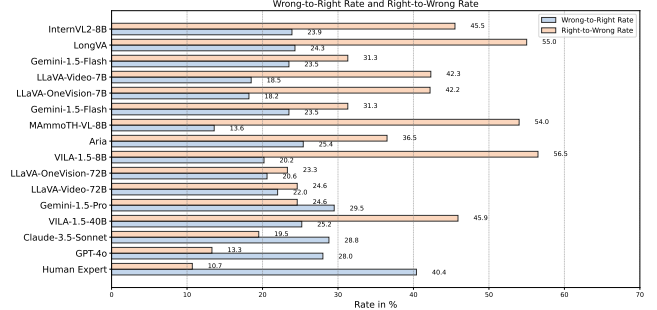


Figure 4. Performance comparison across tracks before and after adding audio transcripts.

As shown in Fig. 4, audio transcripts yield overall performance improvements across different evaluation tracks. In the Comprehension track, the enhancement is most pronounced, reflecting audio’s contribution to video content understanding. Similarly, the Perception track demonstrates



(a) Comparison of  $\Delta_{\text{knowledge}}$  (performance improvement in the Adaptation track after watching the video compared to before).



(b) Comparison of Wrong-to-Right Rate (the percentage of Adaptation track questions that were initially answered incorrectly without the video but correctly after watching the video) and Right-to-Wrong Rate (vice versa).

Figure 5. Key findings in the experiment of  $\Delta_{\text{knowledge}}$ .

performance gains, suggesting that audio enhances information extraction from videos. The Adaptation track, however, shows a decrease in performance. This decline indicates that while audio enriches basic understanding, it might complicate the adaptation of knowledge to novel scenarios. These contrasting effects reveal a trade-off: audio transcripts enhance immediate comprehension but potentially constrain models’ ability to adapt knowledge to new scenarios.

## 5. Knowledge Acquisition in Adaptation Track

### 5.1. Settings

We introduce a knowledge acquisition metric  $\Delta_{\text{knowledge}}$  to measure how much knowledge models gain from videos through their performance improvement on practical exam questions in the Adaptation track. We define  $\Delta_{\text{knowledge}}$  as:

$$\Delta_{\text{knowledge}} = \frac{\text{Acc}_{\text{post}} - \text{Acc}_{\text{pre}}}{100\% - \text{Acc}_{\text{pre}}} \times 100\%$$

where  $\text{Acc}_{\text{pre}}$  and  $\text{Acc}_{\text{post}}$  represent the accuracy before and after watching the video, respectively. This normalized metric accounts for different baseline difficulty levels. For example, improving from 90% to 95% ( $\Delta_{\text{knowledge}} = 50\%$ ) indicates more substantial video-based learning than improving from 0% to 5% ( $\Delta_{\text{knowledge}} = 5\%$ ). We evaluate  $\Delta_{\text{knowledge}}$  across top-performing open-source and proprietary models.

### 5.2. Findings

**Human-Model Knowledge Acquisition Gap:** Fig. 5a reveals a substantial disparity between human and model learning capabilities. Humans demonstrate a  $\Delta_{\text{knowledge}}$  of 33.1% after viewing the videos, while the best-performing model GPT-4o achieves only 15.6%. Some models even exhibit negative  $\Delta_{\text{knowledge}}$ , suggesting their performance declines after video exposure.

This gap highlights a fundamental challenge in current models. Humans naturally acquire information through video-based learning, having developed this capability through classroom education and video content throughout

their lives. While many models can process video information, they struggle to effectively learn new knowledge from the video and apply it in practice.

**Video Impact on Model Responses:** While low  $\Delta_{\text{knowledge}}$  scores might suggest limited net knowledge gain, models’ responses change substantially after watching the videos. As shown in Fig. 5b, we analyze these changes through two metrics: Wrong-to-Right Rate (the percentage of questions initially answered incorrectly but correctly after watching videos) and Right-to-Wrong Rate (the percentage of questions correctly answered before but incorrectly after watching videos). We define the Wrong-to-Right Rate as:

$$\text{Wrong-to-Right Rate} = \frac{N_{\text{Wrong-to-Right}}}{N_{\text{Wrong-before}}} \times 100\%$$

, where  $N_{\text{Wrong-to-Right}}$  refers to the number of questions that were answered incorrectly before watching the video but correctly after watching the video, and  $N_{\text{Wrong-before}}$  is the total number of questions that were answered incorrectly before watching the video.

Similarly, we define the Right-to-Wrong Rate as:

$$\text{Right-to-Wrong Rate} = \frac{N_{\text{Right-to-Wrong}}}{N_{\text{Right-before}}} \times 100\%$$

, where  $N_{\text{Right-to-Wrong}}$  refers to the number of questions that were answered correctly before watching the video but incorrectly after watching the video, and  $N_{\text{Right-before}}$  is the total number of questions that were answered correctly before watching the video.

Interestingly, models achieve moderate Wrong-to-Right Rates (e.g., Gemini-1.5-Pro: 29.5%), demonstrating certain ability to acquire knowledge from videos. However, their high Right-to-Wrong Rates (e.g., LongVA: 55.0%) significantly offset these gains, indicating that they struggle to maintain their initial correct answers while processing new video information. In contrast, human experts demon-

✘ **Method Adaptation Error**

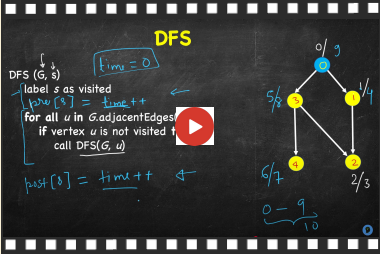
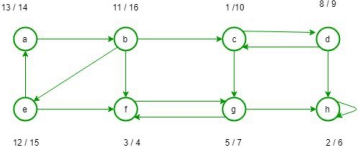
Adaptation Track Question "Engineering – Computer Science"	Video Lecture "Pre and Post visited Times in DFS"
<p><b>Question:</b> In <b>Depth First Search (DFS)</b>, discovery/finishing time stamps are <math>x/y</math>, where <math>x</math> is discovery time and <math>y</math> is finishing time. If the time stamps of <b>g</b> and <b>h</b> in the graph are changed to <b>2/7</b> and <b>5/6</b>, respectively, which of the following traversals are valid depth first forests in &lt;image 1&gt;?</p> <p><b>Options:</b></p> <p>(A) {a, b, c, d, e} {f, g, h}    (B) {a, b, e} {f, g} {c, d} {h}</p> <p>(C) {a, b, e, h} {c, d, f, g}    (D) {a, b, e} {c, d, h} {f, g}</p> <p>(E) {a, b, e} {c, d, h} {f, g}    (F) <b>{a, b, e} {c, d, f, g, h}</b> ✓</p> <p>(G) {a, b, e, f, g} {c, d, h}    (H) {a, b, e} {c, g, f, h} {d}</p> <p>(I) {a, b, c, e} {d, f, g, h}    (J) {a, b, f, e} {c, g, h, d}</p>	
 <p>&lt;image 1&gt;</p>	
Before Video	After Video
<p><b>Claude-3.5-Sonnet:</b></p> <p>(1) Given the new timestamps and analyzing the connectivity and timing constraints: <b>h must be in a separate tree from g</b> due to their new timestamps and the cycles, c and d form a strongly connected component, a, b, and e form another connected component, f and g are connected.</p> <p>(2) The answer is D: {a, b, e} {c, d, h} {f, g}. This is correct because... ✘</p> <p><b>Human:</b></p> <p>(1) Looking at the timestamps, c (1/10) is visited first, followed by g, f, h.</p> <p>(2) Node h has a cycle and does not have any outgoing edges, therefore it must be the end of the first DFS tree. Then <b>d (8/9) was visited and will form a second DFS tree</b>. Then b was discovered followed by e and a.</p> <p>(3) Hence the answer should be H: {a, b, e} {c, g, f, h} {d}. ✘</p>	<p><b>Claude-3.5-Sonnet:</b></p> <p>Let me help solve this step by step using what was taught in the video...</p> <p>(1) First, let's recall the <b>key principles</b>: Vertices in the same DFS tree have overlapping intervals. Parent's interval must fully contain child's interval...</p> <p>(2) Analyzing the graph: - Node a is discovered first in the graph - Nodes a, b, e form one tree component - Nodes c, d form another component... ✘</p> <p>(3) Therefore, the answer is B: {a, b, e} {f, g} {c, d} {h}. ✘</p> <p><b>Human:</b></p> <p>(1) The video helped me understand the <b>principle</b>: the DFS tree's root node's interval (discovery/finish timestamps) should contain all its descendants' intervals, as the root is discovered first and finished in the end.</p> <p>(2) Applying this principle, I could see that c's interval (1/10) contains the intervals of nodes g, f, h, and d, making them all part of c's DFS tree. ✓</p> <p>(3) Hence the answer is F: {a, b, e} {c, d, f, g, h}. ✓</p>
<p><b>Comment:</b></p> <ol style="list-style-type: none"> <li>The video introduces basic DFS principles with timestamps, while the <b>adaptation question tests these concepts in a more complex graph with cycles</b>.</li> <li>Before the video, both Claude and Human incorrectly focused on cycles and connectivity patterns.</li> <li>After learning the video, both started correctly by understanding the core principle: a root node's interval must contain its descendants' intervals. However, <b>Claude demonstrated a Method Adaptation Error</b> - despite understanding the principle, it failed to correctly <b>adapt it to the complex graph</b> (as shown in blue). In contrast, <b>Human successfully applied the method</b> of interval containment in this new scenario (as highlighted in the box).</li> <li>This reveals the challenges of adapting the method from the video in novel, real-world scenarios.</li> </ol>	

Figure 6. A Case of Method Adaptation Error. The model can recall the correct knowledge from the video but fails to adapt the method to a new scenario. More error cases are analyzed in the Appendix.

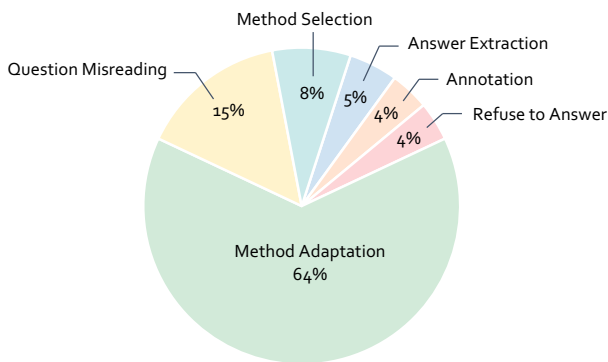


Figure 7. Distribution of the 100 human-annotated errors in Claude-3.5-Sonnet.

strate effective knowledge acquisition with a higher Wrong-to-Right Rate (40.4%) and a lower Right-to-Wrong Rate (10.7%). This indicates humans' ability to integrate new

knowledge while preserving their prior knowledge. These findings highlight a gap between human and model capabilities in video-based learning, particularly in maintaining existing knowledge while processing new information from videos.

### 5.3. Error Analysis

We analyzed the Claude-3.5-Sonnet errors in the Adaptation track by examining 100 randomly sampled error cases. Human annotators analyzed these cases to identify the root causes of mispredictions. The distribution of these errors is shown in Fig. 7, with more error cases provided in the Appendix.

**Method Selection Error (8%):** The model's initial thinking direction is incorrect. For example, the model fails to adopt the correct solution strategy demonstrated in the video.

**Method Adaptation Error (64%):** These represent cases where the model can correctly recall and understand the methods demonstrated in the video but fails to adapt the



method to new scenarios properly. For example, Fig. 6 shows how models can struggle with subtle scenario differences between the video example and the Adaptation question. While the model recalls the core DFS principles from a simple tree example in the video, it fails to adapt these principles flexibly when working with a more complex graph containing cycles. This type of error reveals its limitations in video-based learning when applying the learned methods across different contexts.

**Question Misreading Error (15%):** These errors stem from misinterpreting the question requirements, such as misreading numerical values or conditions. Such errors are unrelated to the model’s ability to apply knowledge from videos.

**Other Errors:** These include Refuse to Answer (4%), where models express uncertainty and decline to provide an answer; Annotation error (4%), where the annotation is inaccurate; and Answer Extraction error (5%), where we failed to extract the answer from the longer output.

Our experiment on  $\Delta_{\text{knowledge}}$  provides insights for future research in knowledge acquisition from videos: **1)** Models showcase certain ability to acquire knowledge from videos, as indicated by their modest Wrong-to-Right Rates. However, the high Right-to-Wrong Rates often negate these gains, suggesting that models struggle to retain their initial correct reasoning when processing new information from video. **2)** The Question Misreading and Method Selection Errors highlight the fundamental limitations in processing knowledge-intensive videos. Accurate question interpretation and a thorough understanding of video knowledge are crucial for successful knowledge application. **3)** The significant proportion of Method Adaptation errors reveals a gap between comprehension and adaptation capabilities, suggesting that applying the knowledge from videos to solve a novel, practical scenario remains challenging for the current models.

## 6. Conclusion

Video-MMMU systematically evaluates how large multi-modal models (LMMs) acquire knowledge from videos through three cognitive stages: Perception, Comprehension, and Adaptation. Through our proposed  $\Delta_{\text{knowledge}}$  metric, we reveal a gap between human and model performance, particularly in adapting acquired knowledge to novel, practical scenarios. Our insights from Video-MMMU underscore the critical need for future research to enhance LMMs’ ability to learn and apply knowledge from video content effectively.

## References

- [1] Anthropic. Claude Team. Introducing Claude 3.5 Sonnet. <https://www.anthropic.com/claude/sonnet>, 2024. 2, 5, 6, 1
- [2] Michael Boratko, Harshit Padigela, Divyendra Mikkilineni, Prithish Yuvraj, Rajarshi Das, Andrew McCallum, Maria Chang, Achille Fokoue-Nkoutche, Pavan Kapanipathi, Nicholas Mattei, Ryan Musa, Kartik Talamadupula, and Michael Witbrock. A systematic classification of knowledge, reasoning, and context within the ARC dataset. In *Proceedings of the Workshop on Machine Reading for Question Answering*, pages 60–70, Melbourne, Australia, 2018. Association for Computational Linguistics. 2
- [3] Mu Cai, Reuben Tan, Jianrui Zhang, Bocheng Zou, Kai Zhang, Feng Yao, Fangrui Zhu, Jing Gu, Yiwu Zhong, Yuzhang Shang, Yao Dou, Jaden Park, Jianfeng Gao, Yong Jae Lee, and Jianwei Yang. Temporalbench: Towards fine-grained temporal understanding for multimodal video models. *arXiv preprint arXiv:2410.10818*, 2024. 2
- [4] Wenhao Chai, Enxin Song, Yilun Du, Chenlin Meng, Vashisht Madhavan, Omer Bar-Tal, Jeng-Neng Hwang, Saining Xie, and Christopher D. Manning. Auroracap: Efficient, performant video detailed captioning and a new benchmark. *arXiv preprint arXiv:2410.03051*, 2024. 2
- [5] Xiuyuan Chen, Yuan Lin, Yuchen Zhang, and Weiran Huang. Autoeval-video: An automatic benchmark for assessing large vision language models in open-ended video question answering. *arXiv preprint arXiv:2311.14906*, 2023. 5
- [6] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 5, 6, 2
- [7] Xinyu Fang, Kangrui Mao, Haodong Duan, Xiangyu Zhao, Yining Li, Dahua Lin, and Kai Chen. Mmbench-video: A long-form multi-shot benchmark for holistic video understanding. *arXiv preprint arXiv:2406.14515*, 2024. 2, 5
- [8] Mary Forehand. Bloom’s taxonomy. *Emerging perspectives on learning, teaching, and technology*, 41(4):47–56, 2010. 2
- [9] Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024. 2, 5
- [10] Noa Garcia, Mayu Otani, Chenhui Chu, and Yuta Nakashima. Knowit vqa: Answering knowledge-based questions about videos. In *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020. 2
- [11] Michail N Giannakos. Exploring the video-based learning research: A review of the literature. *British Journal of Educational Technology*, 44(6):E191–E195, 2013. 2
- [12] Jarvis Guo, Tuney Zheng, Yuelin Bai, Bo Li, Yubo Wang, King Zhu, Yizhi Li, Graham Neubig, Wenhu Chen, and Xiang Yue. Mammoth-vl: Eliciting multimodal reasoning with instruction tuning at scale. 2024. 5, 6, 2
- [13] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. 2
- [14] Muhammad Uzair khattak, Muhammad Ferjad Naeem, Jameel Hassan, Naseer Muzzamal, Federcio Tombari, Fahad Shahbaz Khan, and Salman Khan. How good is my video lmm? complex video reasoning and robustness evaluation suite for video-lmms. *arXiv:2405.03690*, 2024. 2

- [15] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 5, 6, 2
- [16] Dongxu Li, Yudong Liu, Haoning Wu, Yue Wang, Zhiqi Shen, Bowen Qu, Xinyao Niu, Guoyin Wang, Bei Chen, and Junnan Li. Aria: An open multimodal native mixture-of-experts model. *arXiv preprint arXiv:2410.05993*, 2024. 5, 6, 2
- [17] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, Limin Wang, and Yu Qiao. Mvbench: A comprehensive multimodal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22195–22206, 2024. 5
- [18] Shicheng Li, Lei Li, Shuhuai Ren, Yuanxin Liu, Yi Liu, Rundong Gao, Xu Sun, and Lu Hou. Vitatecs: A diagnostic dataset for temporal concept understanding of video-language models. *arXiv preprint arXiv:2311.17404*, 2024. 2
- [19] Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models. *arXiv preprint arXiv:2312.07533*, 2023. 5, 6, 2
- [20] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv: 2403.00476*, 2024. 2, 5
- [21] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [22] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. In *Advances in Neural Information Processing Systems*, pages 46212–46244. Curran Associates, Inc., 2023. 2
- [23] Meta. Llama 3.2: Revolutionizing Edge AI and Vision with Open, Customizable Models. <https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/>, 2024. 5, 6
- [24] Meredith Ringel Morris, Jascha Sohl-Dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksandra Faust, Clement Farabet, and Shane Legg. Position: Levels of AGI for operationalizing progress on the path to AGI. In *Proceedings of the 41st International Conference on Machine Learning*, pages 36308–36321. PMLR, 2024. 2
- [25] Munan Ning, Bin Zhu, Yujia Xie, Bin Lin, Jiayi Cui, Lu Yuan, Dongdong Chen, and Li Yuan. Video-bench: A comprehensive benchmark and toolkit for evaluating video-based large language models. *arXiv preprint arXiv:2311.16103*, 2023. 5
- [26] OpenAI. Introducing openai o1. <https://openai.com/o1/>, 2024. 4
- [27] OpenAI. Hello gpt4-o. <https://openai.com/index/hello-gpt-4o/>, 2024. 2, 3, 5, 6, 1
- [28] Viorica Pătrăucean, Lucas Smaira, Ankush Gupta, Adrià Recasens Contente, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Joseph Heyward, Mateusz Malinowski, Yi Yang, Carl Doersch, Tatiana Matejovicova, Yury Sulsky, Antoine Miech, Alex Frechette, Hanna Klimczak, Raphael Koster, Junlin Zhang, Stephanie Winkler, Yusuf Aytar, Simon Osindero, Dima Damen, Andrew Zisserman, and João Carreira. Perception test: A diagnostic benchmark for multimodal video models. In *Advances in Neural Information Processing Systems*, 2023. 2
- [29] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*, 2022. 6
- [30] Marija Sablić, Ana Miroslavljević, and Alma Škugor. Video-based learning (vbl)—past, present and future: An overview of the research published from 2008 to 2019. *Technology, Knowledge and Learning*, 26(4):1061–1077, 2021. 2
- [31] Enxin Song, Wenhao Chai, Guan hong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, et al. Moviechat: From dense token to sparse memory for long video understanding. *arXiv preprint arXiv:2307.16449*, 2023. 2
- [32] Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 4, 5, 6, 1, 2
- [33] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 5
- [34] Weihang Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Shiyu Huang, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*, 2024. 2
- [35] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 2
- [36] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhramil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024. 2
- [37] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. *arXiv preprint arXiv:2407.15754*, 2024. 2
- [38] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9777–9786, 2021. 2
- [39] Binzhu Xie, Sicheng Zhang, Zitang Zhou, Bo Li, Yuanhan Zhang, Jack Hessel, Jingkang Yang, and Ziwei Liu. Funqa: Towards surprising video comprehension. In *European Con-*

- ference on Computer Vision*, pages 39–57. Springer, 2025. [2](#)
- [40] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *ACM Multimedia*, 2017. [2](#)
- [41] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [2](#)
- [42] Kexin Yi\*, Chuang Gan\*, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B. Tenenbaum. Clevrer: Collision events for video representation and reasoning. In *International Conference on Learning Representations*, 2020. [2](#)
- [43] Ahmed Mohamed Fahmy Yousef, Mohamed Amine Chatti, and Ulrik Schroeder. The state of video-based learning: A review and future perspectives. *International Journal on Advances in Life Sciences*, 6(3):122–135, 2014. [2](#)
- [44] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *AAAI*, pages 9127–9134, 2019. [2](#)
- [45] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*, 2024. [2](#), [4](#)
- [46] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhua Chen, and Graham Neubig. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. *arXiv preprint arXiv:2409.02813*, 2024. [2](#), [4](#)
- [47] Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Reality check on the evaluation of large multimodal models. *arXiv preprint arXiv:2407.12772*, 2024. [5](#)
- [48] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision. *arXiv preprint arXiv:2406.16852*, 2024. [5](#), [6](#), [2](#)
- [49] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. [5](#), [6](#), [2](#)
- [50] Yuanhan Zhang, Kaichen Zhang, Bo Li, Fanyi Pu, Christopher Arif Setiadharm, Jingkang Yang, and Ziwei Liu. Worldqa: Multimodal world knowledge in videos through long-chain reasoning. *arXiv preprint arXiv:2405.03272*, 2024. [2](#)
- [51] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. AGIEval: A human-centric benchmark for evaluating foundation models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2299–2314, Mexico City, Mexico, 2024. Association for Computational Linguistics. [2](#)
- [52] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video understanding. *arXiv preprint arXiv:2406.04264*, 2024. [2](#)
- [53] Luowei Zhou, Chenliang Xu, and Jason J Corso. Towards automatic learning of procedures from web instructional videos. In *AAAI Conference on Artificial Intelligence*, pages 7590–7598, 2018. [2](#)

# Video-MMMU: Evaluating Knowledge Acquisition from Multi-Discipline Professional Videos

## Supplementary Material

### 7. Subjects by Discipline

Discipline	Subjects
Art	Art History, Art Theory, Design, Music
Business	Accounting, Economics, Finance, Manage, Marketing
Science	Biology, Chemistry, Geography, Math, Physics
Medicine	Basic Medical Science, Clinical Medicine, Diagnostics and Laboratory Medicine, Pharmacy, Public Health
Humanities	History, Literature, Psychology, Sociology
Engineering	Agriculture, Architecture and Engineering, Computer Science, Electronics, Energy and Power, Materials, Mechanical Engineering

Table 3. Subjects categorized under six disciplines.

### 8. Additional Knowledge Acquisition Experiment Results

We present the results of the  $\Delta_{\text{knowledge}}$  experiment in Table 4. This table includes a detailed breakdown of the number of questions that transitioned from Wrong-to-Right and Right-to-Wrong, along with the corresponding rates.

The  $\Delta_{\text{knowledge}}$  metric reveals a gap between human experts and models, particularly in their ability to learn new information from videos. This skill, which humans exhibit naturally through video-based learning, arises from our long-standing reliance on videos as a medium for acquiring knowledge. Humans have developed a proficiency in extracting, retaining, and applying information from visual content, making video learning an essential component of natural knowledge acquisition.

For models to perform effectively in real-world environments, the ability to learn and adapt from videos is crucial. This capability would allow models to continuously evolve and refine their understanding, thereby enhancing their utility in dynamic and complex scenarios. However, the result suggests that current models are not yet capable of effectively acquiring new knowledge from video and applying it in practice. This suggests that future research needs to focus on improving how models acquire knowledge from videos -

#### System Message:

As an AI assistant, you should watch and learn from the video. Then, adapt what you learned to answer the following question. The image for this question is at the end of the video.

Question: [Question Text]

Options:

A) [Option A]

B) [Option B]

[etc.]

Figure 8. Prompt for Adaptation track.

specifically, their ability to understand, remember, and apply information from video content. These improvements will be crucial for future LMMs to work effectively in the wild.

### 9. Prompt for Adaptation Track

In the adaptation track, we append the question’s image to the end of each video. We introduce the prompt as shown in Fig. 8.

### 10. Prompt for Determining the Helpfulness of Audio

For all samples in Video-MMMU, we employ Gemini 1.5 Pro [32] to analyze each video-question pair and determine if audio might be helpful to solve the question, as shown in Fig. 3c. This analysis will benefit more future Large Multi-modal Models (LMMs) with audio processing capabilities. We introduce the prompt as shown in Fig. 9.

### 11. Annotation Pipeline

We illustrate our pipeline for video collection and QA annotation in Fig. 10.

### 12. More Error Analysis

This section presents a comprehensive analysis of error cases across all three tracks. We begin by examining errors made by Claude-3.5-Sonnet [1] in the Adaptation track. Specifically, Fig. 11 illustrates Method Selection Errors, while Fig. 12 demonstrates Question Misreading Errors.

We also analyze error cases by GPT-4o [27] in the Adaptation track. Fig. 13 and Fig. 14 present Method Adaptation Error and Question Misreading Error, respectively.

Model	$\Delta_{\text{knowledge}}$ (%)	Wrong-to-Right		Right-to-Wrong	
		No. of Questions	Rate (%)	No. of Questions	Rate (%)
Human Expert	33.1	72	40.4	13	10.7
GPT-4o [27]	15.6	44	28.0	19	13.3
Claude-3.5-Sonnet [1]	11.4	42	28.8	30	19.5
VILA-1.5-40B [19]	9.4	57	25.2	34	45.9
Gemini-1.5-Pro [32]	8.7	49	29.5	33	24.6
LLaVA-Video-72B [49]	7.1	40	22.0	29	24.6
LLaVA-OneVision-72B [15]	6.6	37	20.6	28	23.3
VILA-1.5-8B [19]	5.9	48	20.2	35	56.5
Aria [16]	3.2	47	25.4	42	36.5
MAmmoTH-VL-8B [12]	1.5	48	23.9	45	45.5
Gemini-1.5-Flash [32]	-3.3	39	23.5	42	31.3
LLaVA-Video-7B [49]	-5.3	35	18.5	47	42.3
LLaVA-OneVision-7B [15]	-5.6	36	18.2	43	42.2
LongVA [48]	-7.0	29	13.6	47	54.0
InternVL2-8B [6]	-8.5	46	24.3	61	55.0

Table 4. Additional Knowledge Acquisition Experiment Results with Delta (%) values.

Furthermore, we investigate error cases in both the Perception and Comprehension tracks. For the Perception track, we present two representative error cases in Fig. 15 and Fig. 16. Similarly, for the Comprehension track, we analyze two error cases shown in Fig. 17 and Fig. 18. Each case study includes a detailed analysis of the observed errors.

### 13. Wrong-to-Right Case Analysis

For the Adaptation track, we also analyze the Wrong-to-Right examples where models successfully learned from video content to correctly solve Adaptation track questions. For Claude-3.5-Sonnet [1], we present three such examples in Fig. 19, Fig. 20, and Fig. 21. Additionally, we present a Wrong-to-Right example of GPT-4o [27] in Fig. 22. Each case study provides a detailed analysis of how the model successfully adapted its knowledge.

**System Message:**

```
template = """"\n[System]
```

You are an assistant that helps with question evaluation. I will provide you with a video, along with a pair of questions and answers. Your task is to assess whether the question requires audio information from the video to be answered, or if it can be answered purely through visual information. You need to provide a complete and detailed reason explaining why the question does or does not require audio from the video.

```
[Question]\n{question}
```

```
[Answer]\n{answer}
```

```
[Standard]
```

The standard for determining whether audio is necessary is: if a question does not require audio, then I should be able to turn off the video's sound and still be able to infer the correct answer entirely from the visual information.

```
[Output Format]
```

Your answer must strictly follow the JSON format below:

```
{\n  "reason": "This question requires audio information from the video to be answered because...",\n  "use_audio": true\n}
```

"use\_audio" should be set to "true" if the question requires audio information from the video to be answered, and "false" otherwise.

Please note that you should output only the JSON code, with no additional information.\n"""

Figure 9. Prompt for determining the helpfulness of audio.

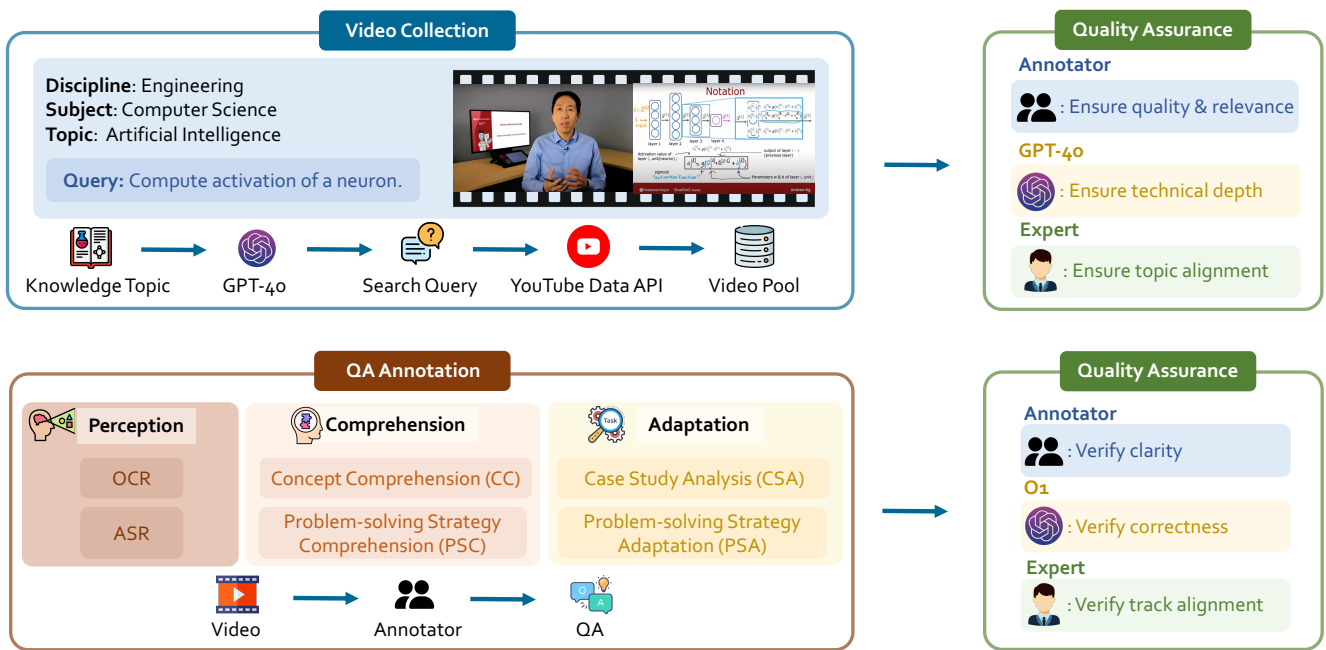


Figure 10. An illustration of the dataset curation pipeline.

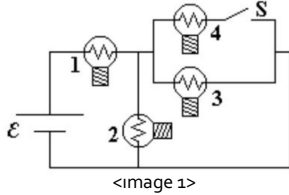


### Method Selection Error

#### Video Lecture --- "Physics Help: Series and Parallel Circuits Electricity Diagrams Part 4"



#### Adaptation Question



**Question:** <image 1>A circuit is connected as shown. All light bulbs are identical. When the switch in the circuit is closed illuminating bulb #4, which other bulb(s) also become brighter?

**Option:**

- (A) Bulb #1 only
- (B) Bulb #2 only
- (C) Bulbs #2 and #3 only
- (D) Bulbs #1, #2, and #3
- (E) Bulbs #1 and #2 only"
- (F) Bulbs #1 and #3 only
- (G) Bulb #4
- (H) None
- (I) All
- (J) Not sure

#### Before Video

##### Claude-3.5-Sonnet:

When the switch is closed, bulb #4 is added in parallel with bulb #3. This changes the total resistance in the circuit and affects the brightness of the other bulbs.

- When the switch is closed, bulbs #3 and #4 are in parallel, and this parallel combination is in series with bulbs #1 and #2.
- Closing the switch reduces the total resistance of the circuit because the parallel combination of bulbs #3 and #4 has a lower resistance than bulb #3 alone. This increases the total current from the battery.
- Since bulbs #1 and #2 are in series with the parallel combination of bulbs #3 and #4, the increased total current flows through bulbs #1 and #2, making them brighter. Therefore, the correct answer is: **F. Bulbs #1 and #2 only.**

#### After Video

##### Claude-3.5-Sonnet:

To determine which bulbs become brighter when the switch is closed, we need to analyze the circuit configuration and the effect of closing the switch on the current distribution. When the switch (S) is closed:

- It creates a parallel path with bulb #3
- The total resistance in that branch decreases (parallel resistors have less total resistance)
- This means more current will flow through the main circuit. Following the circuit:\n- Bulb #1 will get brighter because more total current flows through the main circuit\n- Bulb #2 will get brighter for the same reason as Bulb #1\n- Bulb #3 will actually get dimmer because some of its current is now diverted through Bulb #4
- Therefore, when switch S is closed and Bulb #4 lights up, Bulbs #1 and #2 will become brighter, while Bulb #3 becomes dimmer. The correct answer is **F: Bulbs #1 and #2 only.**

##### Human:

When the switch is closed, Bulb #4 is added in parallel with Bulb #3. This means the current through the circuit increases overall. Since all the bulbs are identical, the increased current is equally shared among all the bulbs in the circuit. Therefore, all the bulbs including Bulbs #1, #2, and #3 will become brighter because the total power delivered to the circuit has increased. Hence the answer is **D. Bulbs #1, #2, and #3.**

##### Human:

As explained in the video, we can calculate the current through each bulb to determine the brightness.

Before switch closes: Total  $R = R + (R \times R)/(R+R) = 3R/2$   
 Initial current  $I = \epsilon/(3R/2) = 2\epsilon/3R$ ,  $I_1 = 2\epsilon/3R$ ,  $I_2 = I_3 = \epsilon/3R$   
 After switch closes: Total  $R = R + (R/2 \times R)/(R/2 + R) = 4R/3$   
 New current  $I = \epsilon/(4R/3) = 3\epsilon/4R$ ,  $I_1 = 3\epsilon/4R$  (increased),  $I_2 = I_3 = \epsilon/4R$  (decreased)  
 $I_4 = \epsilon/4R$  (turns on)  
 Only  $I_1$  increases, hence we choose **A. Bulb #1 only.**

##### Comment:

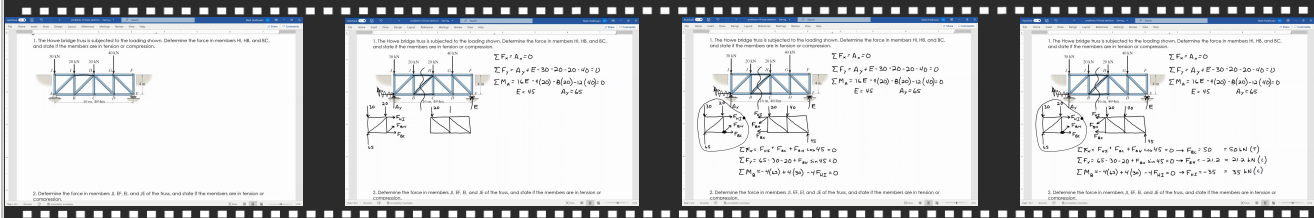
- The video introduces a **problem-solving approach: use Ohm's law to calculate current through each bulb to determine brightness.** The adaptation question presents a new scenario with an added switch, but we can use the identical approach to determine the current through each bulb.
- Although seeing the same circuit problem, the model **failed to adopt the video's quantitative approach of calculating currents.** Instead, it **still used qualitative analysis.** This might be insufficient to determine current change for each bulb, leading to wrong answer.
- In contrast, human effectively adopts the video's approach to the new scenario. Human calculate the actual currents ( $\epsilon/R$ ) through each bulb before/after switch closes, which clearly shows only bulb #1 increases.
- In essence, while the same problem-solving approach (calculate currents) works in both scenarios, the model still uses its own qualitative analysis instead of learning from the video, leading to incorrect conclusions.

Figure 11. A sample error case in the Adaptation track: Method Selection Error by Claude-3.5-Sonnet.

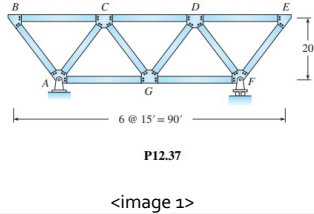


**Question Misreading Error**

**Video Lecture – “Photoelectric Effect Graphs”**



**Adaptation Question**



**Question:** If the truss is to be designed for a uniform live load of 0.32 kip/ft that can be placed anywhere on the span in addition to a concentrated live load of 24 kips that can be positioned where it will produce the largest force in bar CG, determine the maximum value of live load force (tension, compression, or both) created in bar CG. <image 1>

- Option:**  
 (A) 9.75 kips compression (B) 5.75 kips compression (C) 7.75 kips compression  
 (D) 10.25 kips compression (E) 6.25 kips compression (F) 8.25 kips compression  
 (G) 11.50 kips compression (H) 4.50 kips compression (I) 9.00 kips compression (J) 7.00 kips compression

**Before Video**

**Claude-3.5-Sonnet:**

Let me help solve this truss problem step by step.

- 1) First, let's identify what we're looking for:...
- 2) Key information: - Span = 90 feet Height = 20 feet - **Bar CG is a diagonal member** ✓
- 3) Analysis: - For bar CG, the worst case scenario will be when: The uniform load covers the entire span. The 24-kip concentrated load is placed at joint C.
- 4) Calculations: Total uniform load = 0.32 kip/ft × 90 ft = 28.8 kips. Load per panel point = 28.8/6 = 4.8 kips. The 24-kip concentrated load at C will create the maximum compression
- 5) Using method of sections or influence lines, the force in CG can be determined. The answer is **A: 9.75 kips compression**. ✓

**After Video**

**Claude-3.5-Sonnet:**

Let me help analyze this step by step based on the video content and the truss problem shown.

- For bar CG in the truss:
- 1) We need to consider two types of live loads:
    - Uniform live load of 0.32 kip/ft
    - Concentrated live load of 24 kips
  - 2) Key points to consider:
    - **Bar CG is a vertical member** as shown in the image ✗
    - The truss span is 90 feet and the height is 20 feet
  - 3) For maximum force in CG: Need to consider influence lines for the force...
  - 4) Looking at the given options and considering: The method of sections... Based on the video content and analysis, the maximum compressive force in bar CG would be: **C: 7.75 kips compression** ✗

**Human:**

We can first calculate the Total Uniform Load: 28.8 kips  
 Then we calculate the Reactions Due to Uniform Load:  
 1. Since the load is uniformly distributed, the reactions at A and F will be equal.  
 2.  $R_A = R_F = 14.4$  kips  
 To maximize the force in bar CG, place the 24 kip load at point C (mid-span).  
 We can calculate the Reactions Due to Concentrated Load:  
 The concentrated load of 24 kips at point C will be equally distributed to supports A and F...  
 The maximum force in bar CG is **G. 11.50 kips in compression**. ✗

**Human:**

For maximum force in CG:  
 • Concentrated load should be placed at point C  
 • Uniform load should be placed from A to F (full span)  
 • This creates maximum compression in CG ✓  
 Calculations:  
 Total uniform load = 0.32 kip/ft × 90 ft = 28.8 kips  
 Total load = 28.8 kips + 24 kips = 52.8 kips  
 Due to truss geometry and load position, bar CG will experience compression  
 Using method of sections and equilibrium equations, max compression in CG = 9.75 kips.  
 The answer should be **A: 9.75 kips compression**. ✓

**Comment:**

The video explains an example problem about truss analysis and load distribution. The Adaptation track question presented is a similar question to the video example, focusing on finding the maximum force in bar CG under specified loading conditions.

For the model:

- Before Video: Claude correctly perceived CG as a **diagonal member**, it showed proper initial perception without video.
- After Video: Claude adopted correct method from the video, BUT misidentified CG as a **vertical member**. This perceptual error led to incorrect calculations.

In comparison, human maintained correct perception of CG as diagonal. And before the video, human did not have a correct reasoning. After learning from the video, human effectively integrated the video knowledge and arrived at accurate force calculations.

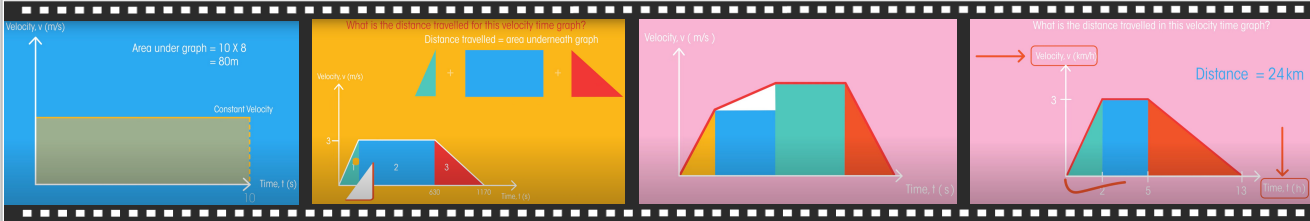
This case reveals an interesting challenge: a visual detail that was correctly perceived in a static image became misinterpreted when presented in the context of video content. **The complexity of processing the long video context** may pose additional challenge to the accurate perception of information.

Figure 12. A sample error case in the Adaptation track: Question Misreading Error by Claude-3.5-Sonnet.

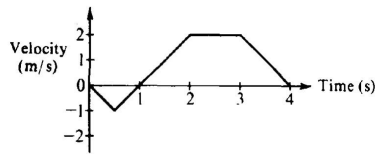


### Method Adaptation Error

#### Video Lecture --- "Area Under Velocity Time Graphs | Forces & Motion"



#### Adaptation Question



<image 1>

**Question:** <image 1>The graph above shows the velocity versus time for an object moving in a straight line. At what time after  $t = 0$  does the object again pass through its initial position?

**Option:**

- (A) 1 s
- (B) Between 1 and 2 s ✓
- (C) 2 s
- (D) Between 2 and 3 s
- (E) 3 s
- (F) Between 3 and 4 s
- (G) 4 s
- (H) Between 4 and 5 s
- (I) 5 s
- (J) Between 5 and 6 s

#### Before Video

##### GPT-4o:

From  $t = 0$  to  $t = 1$  second, the area under the curve is a triangle with a base of 1 second and a height of -1 m/s. The area (displacement) is -0.5.  
 From  $t = 1$  to  $t = 2$  seconds, the area under the curve is a triangle with a base of 1 second and a height of 2 m/s. The area (displacement) is: 1  
 From  $t = 2$  to  $t = 3$  seconds, the velocity is constant at 2 m/s, and the area under the curve is a rectangle with a base of 1 second and a height of 2 m/s. The area (displacement) is: 2.  
 .....

The object does not pass through its initial position within the given time frame of 4 seconds. However, it is possible that it passes through the initial position after 4 seconds. Hence, the closest answer would be **(H) Between 4 and 5 s**. But since the graph does not show this, we cannot definitively determine the exact time from the provided graph. ✗

##### Human:

As shown in the graph, the position should be the point where all the area under the curve become zero. Hence the answer is **(B) Between 1 and 2 s**. ✓

##### Comment:

- The video introduces the concept of calculating displacement using the area under a velocity-time graph. The adaptation question asks for the point where the object returns to its initial position, which means the net displacement must equal zero.
- The examples in the video only show examples with curves above zero, while this question involves a scenario with both positive and negative areas.
- The model demonstrates an understanding of the method by calculating the area under each interval but fails to correctly apply it during reasoning, mistakenly concluding that the object cannot return to its initial position. In contrast, the human accurately adapts the knowledge from the video to this new scenario, correctly summing up the positive and negative areas to determine that the object returns to its initial position between 1 and 2 seconds, showcasing a concise and correct reasoning process.

#### After Video

##### GPT-4o:

To determine when the object passes through its initial position again, we need to consider the area under the velocity-time graph, which represents the displacement of the object. The object will pass through its initial position when the net displacement is zero. ✓

From 0 to 1 second: The object has a negative velocity, so the area is below the time axis. Area =  $0.5 * \text{base} * \text{height} = 0.5 * 1 * (-1) = -0.5 \text{ m}$ .  
 From 1 to 2 seconds: The area is a triangle with base 1 second and height 2 m/s. Area =  $0.5 * \text{base} * \text{height} = 0.5 * 1 * 2 = 1 \text{ m}$ .  
 .....

As all the area is above zero, we cannot derive a position when the net displacement is zero. The object does **not** return to its initial position within the 4 seconds shown in the graph. Based on the given options, the closest answer would be: **J. Between 5 and 6 seconds**. ✗

##### Human:

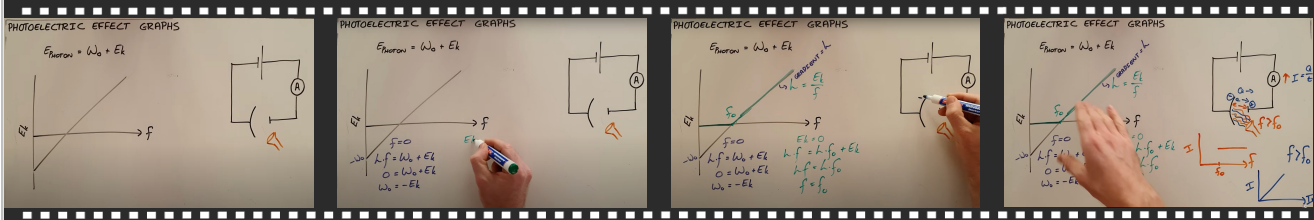
As explained in the video, the area under the curve is the displacement of the object. Hence, we can find the position where the area under the curve is zero. From 0 to 1 seconds, the area is -0.5. From 1 to 2 seconds, the area is 1. Hence the position is between 1 and 2 seconds. The answer should be **(B) Between 1 and 2 s**. ✓

Figure 13. A sample error case in the Adaptation track: Method Adaptation Error by GPT-4o.

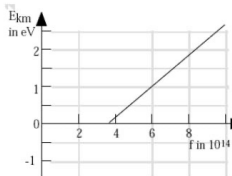


### Question Misreading Error

#### Video Lecture – "Photoelectric Effect Graphs"



#### Adaptation Question



<image 1>

**Question:** <image 1>A student performs the photoelectric effect experiment and obtains the data depicted in the accompanying graph of  $E_{k\max}$  (maximum kinetic energy) of photoelectrons v. the frequency of the photons. What is the approximate work function of this material?

**Option:**

- (A) 1.5 eV ✓
- (D) 3.5 eV
- (G) 3.0 eV

- (B) 2.0 eV
- (E) 1.8 eV
- (H) 3.8 eV

- (C) 2.7 eV
- (F) 2.2 eV
- (I) 2.5 eV

- (J) 4.0 eV

#### Before Video

##### GPT-4o:

The equation governing the photoelectric effect is:

$$E_k = hf - \phi \quad \checkmark$$

where:

- $E_k$  is the maximum kinetic energy of the photoelectrons,
- $h$  is Planck's constant,  $h = 4.14 \times 10^{-15} \text{ eV}$ ,
- $f$  is the frequency of the photons,
- $\phi$  is the work function of the material.

The work function  $\phi$  corresponds to the threshold frequency  $f_0$  where the kinetic energy  $E_k = 0$ .

From the graph: The threshold frequency  $f_0 = 4 \times 10^{14} \text{ Hz}$ . ✗

The work function is calculated as:  $\phi = h * f_0$

Substitute  $f_0 \approx 4 \times 10^{14} \text{ Hz}$  and  $h = 4.14 \times 10^{-15} \text{ eV}$ :  $\phi \approx 1.8 \text{ eV}$  ✗

##### Human:

We can derive from the photoelectric effect equation:

$$E_k = hf - \phi$$

To calculate the work function  $\phi$ , the  $E_k$  must be  $E_k = 0$ .

Hence, we can calculate  $\phi$  by substituting  $f_0 = 4 \times 10^{14} \text{ Hz}$  and  $h = 4.14 \times 10^{-15} \text{ eV}$ :  $\phi \approx 1.8 \text{ eV}$  ✗

#### After Video

##### GPT-4o:

To determine the work function  $\phi$  of the material from the graph of maximum kinetic energy  $E_k$  of photoelectrons versus the frequency  $f$  of the photons, we can use the photoelectric equation:

$$E_k = hf - \phi \quad \checkmark$$

This equation represents a straight line  $y = mx + c$ , where:

- $y$  is  $E_k$
- $x$  is  $f$
- $m$  is  $h$
- $c$  is  $-\phi$

Looking at the graph, the y-intercept appears to be at approximately -2.0 eV.

Therefore, the work function  $\phi = -c = 2.0 \text{ eV}$  ✗

##### Human:

According to the video lecture, we can directly solve this work function  $\phi$  by directly find the y-intercept of the line. ✓

The y-intercept of the line in this graph is approximately -1.5 eV. Hence, the work function  $\phi = 1.5 \text{ eV}$ . ✓

##### Comment:

The video introduces a similar example related to Photoelectric Effect Graphs, demonstrating that the work function ( $\phi$ ) can be simply calculated by finding the y-intercept of the line. This question can be solved without using formulas, as we only need to recognize that **the y-intercept provides the answer directly**.

Before the video:

- Both humans and model approached the problem by using formulas and substituting values
- This indirect method led to incorrect answers and was not the most efficient approach

After the video:

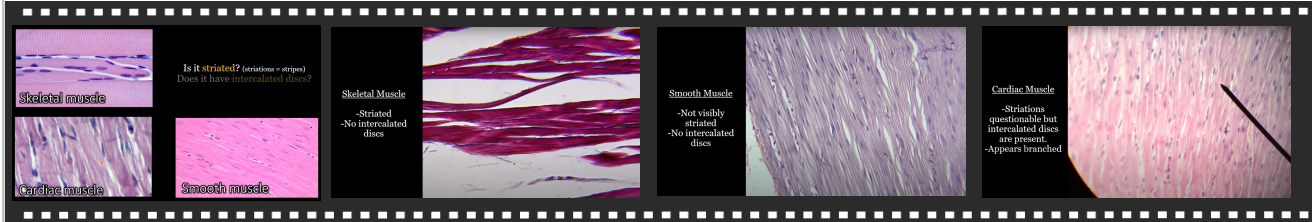
- The model correctly understood that it should look at the y-intercept of the graph.
- However, **GPT-4o misread the graph**, identifying the intercept as -2.0 when it was actually -1.5. This error may have occurred because **GPT-4o incorrectly assumed the x-intercept was at 4**. If the x-intercept were indeed at 4, then a y-intercept of -2.0 might have been plausible
- Humans, on the other hand, demonstrated concise and accurate reasoning. They correctly identified the y-intercept to find the answer of 1.5.

This case represents a Question Misreading Error by GPT-4o. While GPT-4o adopted the correct method and reasoning approach after watching the video, it simply misread the graph.

Figure 14. A sample error case in the Adaptation track: Question Misreading Error by GPT-4o.

**Error in Perception Track**

**Video Lecture – "Identifying Muscle | Review and Practice"**



**Perception Question**

**Question:**

At the beginning of the video, what are the muscles in the lower left corner, upper left corner, and lower right corner, respectively?

**Option:**

- A. Cardiac muscle, Smooth muscle, Skeletal muscle;
- B. Skeletal muscle, Cardiac muscle, Smooth muscle
- C. Skeletal muscle, Smooth muscle, Cardiac muscle;
- D. Smooth muscle, Cardiac muscle, Skeletal muscle
- E. Smooth muscle, Skeletal muscle, Cardiac muscle;
- F. Smooth muscle, Cardiac muscle, Cardiac muscle
- G. Skeletal muscle, Skeletal muscle, Smooth muscle;
- H. Cardiac muscle, Smooth muscle, Smooth muscle
- I. Skeletal muscle, Smooth muscle, Smooth muscle;
- J. Cardiac muscle, Skeletal muscle, Smooth muscle** ✓

**Model response**

**Claude-3.5-Sonnet:**

Looking at the initial frame of the video, the arrangement of the muscle types is: **B. Skeletal muscle (upper left), Cardiac muscle (lower left), Smooth muscle (lower right)**. ✗

**Comment:**

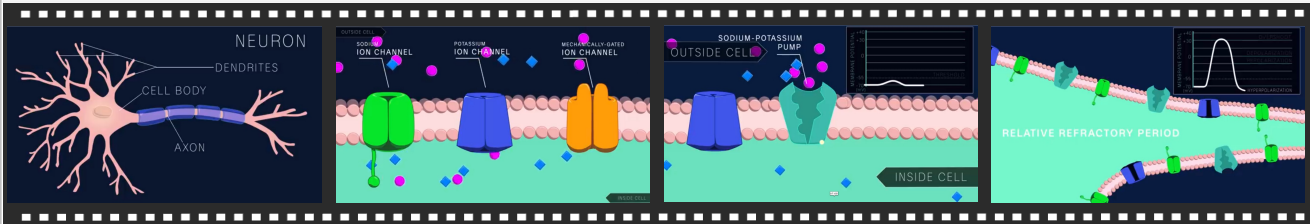
- This example is an error in the perception track.
- The user specifically arranged the **question** to test the model's ability to recognize the **positions of the muscles** when the **sequence is deliberately randomized (lower left, upper left, lower right)**.
- However, the model interpreted the **positions** in a **habitual reading order (left to right, top to bottom)** and provided an **incorrect sequence** of muscle identifications. This highlights a gap in the model's capacity to accurately process and respond to spatial arrangements in visual content.
- Future development should focus on enhancing the model's spatial perception capabilities to prevent such misjudgements.

Figure 15. A sample error case in the Perception track.



### Error in Perception Track

#### Video Lecture – "Action Potential in the Neuron"



#### Perception Question

##### Question:

In the video, the blue ion channel is \_\_\_\_\_, the yellow ion channel is \_\_\_\_\_, and the green ion channel is \_\_\_\_\_.

##### Option:

- A. **POTASSIUM ION CHANNEL, MECHANICALLY-GATED ION CHANNEL, SODIUM ION CHANNEL** ✓
- B. POTASSIUM ION CHANNEL, SODIUM ION CHANNEL, MECHANICALLY-GATED ION CHANNEL
- C. MECHANICALLY-GATED ION CHANNEL, POTASSIUM ION CHANNEL, SODIUM ION CHANNEL
- D. MECHANICALLY-GATED ION CHANNEL, SODIUM ION CHANNEL, POTASSIUM ION CHANNEL
- E. SODIUM ION CHANNEL, POTASSIUM ION CHANNEL, MECHANICALLY-GATED ION CHANNEL
- F. SODIUM ION CHANNEL, MECHANICALLY-GATED ION CHANNEL, POTASSIUM ION CHANNEL

#### Model response

##### Claude-3.5-Sonnet:

Based on the video animation, the ion channels are colored as follows:-

Blue channel: Voltage-gated sodium channel - Orange/yellow channel: Mechanically-gated ion channel - Green channel: Voltage-gated potassium channel.

Therefore, the correct answer would be D: **SODIUM ION CHANNEL, POTASSIUM ION CHANNEL, MECHANICALLY-GATED ION CHANNEL.** ✗

This can be verified by observing how each channel functions in the animation:- The blue channels allow sodium ions (pink circles) to flow in during depolarization. - The green channels allow potassium ions (blue diamonds) to flow out during repolarization. - The orange/yellow channel responds to mechanical stimuli as shown in the membrane.

##### Comment:

- The model's error in this neuroscience question demonstrates an error in perceiving the correct information.
- The model might be **affected by the animation in the video**, and failed to correct perceive the information as shown in the second frame of the video sequence.

Figure 16. A sample error case in the Perception track.



### Error in Comprehension Track

#### Video Lecture – "Discrete Math II - 11.4.2 Spanning Trees - Breadth First Search"

11.4.2 DISCRETE MATH II

Spanning Trees:  
Breadth-First Search

**Breadth-First Search**  
Breadth-first implies we search horizontally before vertically.

1. Choose a "root" vertex for a rooted graph.
2. Add all edges incident to that vertex. The newly-connected vertices are level 1 in the spanning tree.
3. From each vertex in level 1 (in order), add all edges incident to the vertices that are not already included in the spanning tree.
4. Repeat adding levels and edges until all vertices have been visited.

**Breadth-First Practice**

**Breadth-First Practice**

#### Comprehension Question

**Question:**

In the LAST breadth-first search tree example, if the root node is changed to node F, which two nodes will be at the level 3 of the breadth-first spanning tree (root node is level 1)?

**Option:**

- A. C and I
- C. B and C
- E. B and E
- G. C and E
- I. G and H
- B. G and A
- D. D and H
- F. A and I
- H. A and D
- J. H and I

#### Model response

**GPT-4o:**

To determine which two nodes will be at level 3 of the breadth-first spanning tree when the root node is changed to node F, we need to perform a breadth-first search (BFS) starting from node F.

Let's perform the BFS starting from node F:

**Level 1:** Start with node F.

- Current tree: F

**Level 2:** Add all nodes directly connected to F.

- **Nodes connected to F: B, D, G, H** ❌

- Current tree: F | B D G H

**Level 3:** Add all nodes directly connected to nodes in level 2 that are not already in the tree.

- Nodes connected to B: A, C (F is already in the tree)

- Nodes connected to D: A, E (F is already in the tree)

- Nodes connected to G: C (F is already in the tree)

- Nodes connected to H: I (F is already in the tree)

- Current tree: F | B D G H | A C E I

From the tree, we can see that the nodes at level 3 are A, C, E, and I. Therefore, the correct answer is: **"G. C and E"** ❌

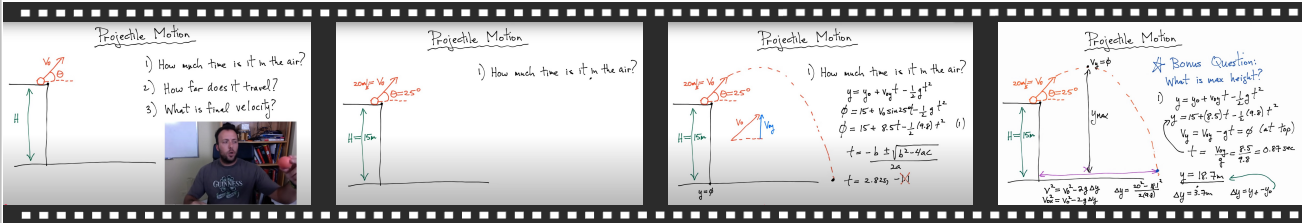
**Comment:**

- The model shows an error in comprehending the problem-solving strategy as demonstrated in the video.
- The video presents two examples of using Breadth-First Search (BFS) to create a spanning tree, which are shown in the third and fourth frames. The question tests understanding by modifying the second example question - specifically by changing the starting node (root) from A to F.
- The model tried to implement BFS but made a mistake in its reasoning. It did not correctly identify which nodes were connected to the new root node F, which led to errors in identifying the level 2 nodes. This mistake then affected the identification of level 3 nodes, resulting in an incorrect answer.
- This error shows that the model struggles to understand and solve problems when the input conditions are changed, even slightly. This limitation reveals a weakness in the model's ability to follow the same reasoning process that was demonstrated in the video. Even though the only change was a different input value, the model should have been able to follow the same process to demonstrated the comprehension of the reasoning process.

Figure 17. A sample error case in the Comprehension track.

**Error in Comprehension Track**

**Video Lecture – "Projectile launched off a cliff at an angle"**



**Comprehension Question**

**Question:**

In the video, Example Question (1) is solved with an angle  $\theta=25$  degrees. If the angle  $\theta$  is adjusted to 30 degrees while all other conditions remain unchanged, what will be the updated result for Example Question (1) as explained in the video?

**Options:**

- (A) 4.00 seconds (B) 2.82 seconds (C) 3.50 seconds (D) 2.50 seconds (E) 3.04 seconds ✓  
 (F) 2.00 seconds (G) 3.15 seconds (H) 1.85 seconds (I) 2.25 seconds (J) 3.85 seconds

**Model response**

**Claude-3.5-Sonnet:**

Let me help you solve this step by step for  $\theta = 30^\circ$ .

- To find the time in the air, we can use the same approach shown in the video, but with  $\theta = 30^\circ$ : Initial velocity ( $V_0$ ) = 20 m/s,  $\theta = 30^\circ$
- First, we need  $V_{0y}$  (initial vertical velocity):  $V_{0y} = V_0 \sin(30^\circ) = 20 \times 0.5 = 10$  m/s
- For time in the air, we can use:
  - **The object returns to its initial height** ✗
  - Vertical velocity at that point equals  $-V_{0y}$
- Using  $V_y = V_{0y} - gt$  ✗
 
$$10 = V_{0y} - gt$$

$$t = 2V_{0y}/g$$

$$t = 2(10)/9.8$$

$$t = 20/9.8$$

$$t = 2.04 \text{ seconds}$$
- Looking at the options provided, the closest answer is: H. 1.85 seconds. Therefore, **H. 1.85 seconds** is the answer. ✗

**Comment:**

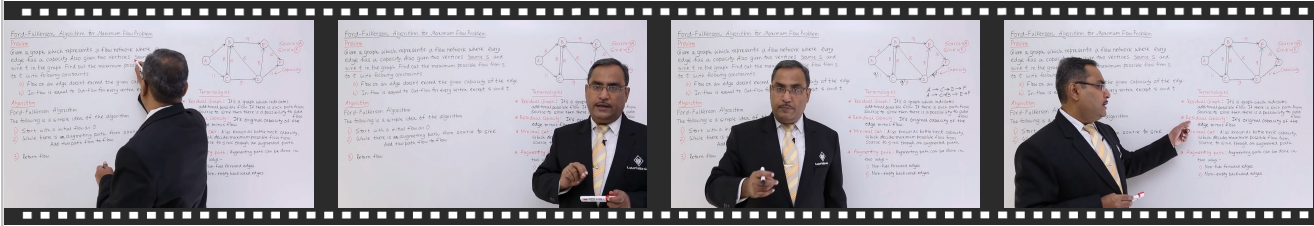
- The case is a comprehension error. This situation clearly demonstrates how the model fails to understand **problem-solving strategies** shown in the video.
- For comprehension questions based on video examples, the test typically uses **the same problem but changes some input values**. To demonstrate true comprehension, the model simply needs to follow the **exact same reasoning process** shown in the video, just with different numbers.
- In this case, the video explains how to calculate a projectile's air time when launched from a cliff at an angle, with the example shown in the second and third frames. The comprehension question uses the identical scenario but only changes one input - the launch angle from 25 to 30 degrees.
- However, the model solved the comprehension question with the **incorrect method and formula**:  $t = (2 * V_{0y}) / g$ . This formula only works for projectiles that return to their starting height, like when launching from ground level. Since this problem involves a 15-meter elevated platform with the projectile landing below its launch point, this shortcut formula does not apply.
- The correct solution** requires using the quadratic equation that accounts for different launch and landing heights:  $y(t) = y_0 + V_{0y} * t - (1/2) * g * t^2$ . This is the same method demonstrated in the video - a process that humans can easily follow but seems challenging for models.
- This example highlights how models struggle to fully comprehend problem-solving strategies shown in educational videos. While a human learner would simply follow the same reasoning steps with new numbers, the model failed to replicate this basic aspect of learning.

Figure 18. A sample error case in the Comprehension track.

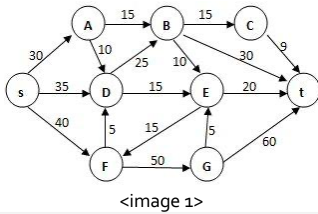


### Wrong-to-Right Example

#### Video Lecture – “Ford Fulkerson Algorithm for Maximum Flow Problem”



#### Adaptation Question



#### Question:

The maximum flow in the network of the given Figure is: <image 1>

#### Option:

- (A) 90
- (B) 98
- (C) 89
- (D) 112
- (E) 101
- (F) 107
- (G) 105
- (H) 119
- (I) 88
- (J) 95 ✓

#### Before Video

##### Claude-3.5-Sonnet:

Let me help you solve this maximum flow problem step by step. To find the maximum flow in this network, I'll use the Ford-Fulkerson algorithm to find augmenting paths from source s to sink t until no more paths are available.

Let's find several paths and their flows:

- 1)  $s \rightarrow A \rightarrow B \rightarrow C \rightarrow t$  ( $\min\{30, 15, 15, 9\} = 9$ )
- 2)  $s \rightarrow D \rightarrow E \rightarrow t$  ( $\min\{35, 15, 20\} = 15$ )
- 3)  $s \rightarrow F \rightarrow G \rightarrow t$  ( $\min\{40, 50, 60\} = 40$ )
- 4)  $s \rightarrow A \rightarrow D \rightarrow E \rightarrow t$  ( $\min\{30, 10, 15, 20\} = 10$ )
- 5)  $s \rightarrow F \rightarrow D \rightarrow E \rightarrow t$  ( $\min\{40, 5, 15, 20\} = 5$ )
- 6)  $s \rightarrow A \rightarrow B \rightarrow E \rightarrow t$  ( $\min\{30, 15, 10, 20\} = 10$ ) ✗
- 7)  $s \rightarrow F \rightarrow G \rightarrow E \rightarrow t$  ( $\min\{40, 50, 5, 20\} = 5$ )
- 8)  $s \rightarrow A \rightarrow B \rightarrow t$  ( $\min\{30, 15, 30\} = 15$ )

Adding up all these flows:  $9 + 15 + 40 + 10 + 5 + 10 + 5 + 15 = 109$   
Looking at the given options, 107 is closest to our calculated value. Therefore, the answer is **F. 107**. ✗

#### After Video

##### Claude-3.5-Sonnet:

Let me help you solve this step by step using the Ford-Fulkerson algorithm that was explained in the video:

To find the maximum flow, we need to:

1. Find augmenting paths from source (S) to sink (t)
2. Calculate the flow through each path based on the minimum capacity along that path
3. Continue until no more augmenting paths are available ✓

Looking at the network diagram, we can identify several possible paths from S to t:  $S \rightarrow A \rightarrow B \rightarrow C \rightarrow t$ ,  $S \rightarrow D \rightarrow E \rightarrow t$ ,  $S \rightarrow F \rightarrow G \rightarrow t$ ,  $S \rightarrow A \rightarrow D \rightarrow E \rightarrow t$ , and ...

After analyzing all possible augmenting paths and calculating the flows while respecting the capacity constraints and flow conservation rules (as explained in the video), the maximum flow that can be achieved in this network is 95 units. Therefore, the answer is **J. 95**. ✓  
This can be verified by checking that:- No path from source to sink can carry more flow. - The flow satisfies capacity constraints on all edges. - Flow conservation is maintained at all vertices except source and sink.

#### Comment:

This example illustrates how the model **successfully learned a problem-solving strategy from a video lecture** on the Ford-Fulkerson Algorithm to correct its initially incorrect answer.

The video demonstrates the proper method for finding augmenting paths and calculating network flows with capacity constraints. The adaptation question tests this knowledge by asking to find the maximum flow in a network diagram.

Before video, the model:

- Double-counted edge capacities and ignored flow constraints
- Added path flows independently, reaching incorrect answer of 107
- Failed to consider residual network capacity

After watching the video, the model:

- Properly identified paths while respecting capacity constraints
- **Applied correct flow calculations** to reach answer of 95 units

In conclusion, this case demonstrates successful knowledge acquisition from video. The model **corrects its understanding of the max flow problems** and the algorithm through learning from the video, and applies the demonstrated algorithm to solve the adaptation problem correctly.

Figure 19. A Wrong-to-Right example of Claude-3.5-Sonnet in the Adaptation track.





## Wrong-to-Right Example

### Video Lecture – “Thin Film Interference Summary”

**Thin Film Interference Summary**

For nearly vertical light rays

air  $n=1.00$   
soap  $n=1.33$   
air  $n=1.00$

No phase change

$2t = m\lambda_{\text{medium}}$  (constructive)  
 $2t = (m + \frac{1}{2})\lambda_{\text{medium}}$  (destructive)

Alternative Format  
 $2t = m\lambda$   
 $2t = (m + \frac{1}{2})\lambda$

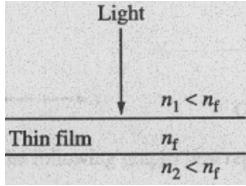
air  $n=1.00$   
soap  $n=1.33$   
glass  $n=1.50$

Phase shift  $\lambda/2$

$2t = m\lambda_{\text{MEDIUM}}$  (destructive)  
 $2t = (m + \frac{1}{2})\lambda_{\text{MEDIUM}}$  (constructive)

Equation	1 phase shift	0 or 2 phase shifts
$2t = (m + \frac{1}{2})\lambda$	constructive	destructive
$2t = m\lambda$	destructive	constructive

### Adaptation Question



<image 1>

**Question:** A thin film with index of refraction  $n_f$  separates two materials, each of which has an index of refraction less than  $n_f$ . A monochromatic beam of light is incident normally on the film, as shown above. If the light has wavelength  $\lambda$  within the film, maximum constructive interference between the incident beam and the reflected beam occurs for which of the following film thicknesses?

**Option:**

- (A)  $2\lambda$  (B)  $\lambda$  (C)  $\lambda/2$  (D)  $\lambda/4$  (E)  $3\lambda$   
 (F)  $5\lambda/2$  (G)  $3\lambda/2$  (H)  $\lambda/8$  (I)  $\lambda/16$  (J)  $4\lambda$

#### Before Video

##### Claude-3.5-Sonnet:

Let me help solve this thin film interference problem.

For constructive interference in a thin film:

- 1) First, we need to consider phase shifts at the boundaries
- 2) Then, consider the path difference for constructive interference

Key points:

- Light is incident normally (perpendicular to surface)  $n_1 < n_f$  and  $n_2 < n_f$  (both surrounding media have lower refractive index than film)
- We want maximum constructive interference
- Phase shifts: At top boundary:  $180^\circ$  ( $\pi$ ) phase shift because  $n_1 < n_f$ . At bottom boundary:  $180^\circ$  ( $\pi$ ) phase shift because  $n_2 < n_f$ . Total phase shift from reflections =  $360^\circ$  ( $2\pi$ ) = **no effective phase shift.**

For constructive interference:

- Total path =  $2t$  (down and up through film)
- Condition:  $2t = m\lambda$  (where  $m$  is an integer)

Therefore,  $t = m\lambda/2$ . The smallest non-zero thickness that satisfies this is:  $t = \lambda/2$ . Looking at the options, **C ( $\lambda/2$ )** is the correct answer.

#### After Video

##### Claude-3.5-Sonnet:

Let me help you solve this step by step using the concepts shown in the video:

1) From the video, we learned that for thin film interference:

- When there's a phase shift at both or neither boundary:  $2t = (m + 1/2)\lambda$  for constructive interference
- When there's a phase shift at only one boundary:  $2t = m\lambda$  for constructive interference

2) In this problem:

- The film has higher index of refraction ( $n_f$ ) than both surrounding materials. Light is incident normally (vertically)

3) Key points:

- When light reflects from a medium with higher  $n$ , there's a phase shift of  $\lambda/2$
- When light reflects from a medium with lower  $n$ , there's no phase shift

4) Analysis:

**Since  $n_f$  is higher than both surrounding materials. This means one phase shift total.**

5) Therefore:

- We use the formula  $2t = (m + 1/2)\lambda$  for constructive interference
- For  $m = 0$ , this gives us  $t = \lambda/4$ . Therefore, the answer is **D:  $\lambda/4$ .**

#### Comment:

This example illustrates how the model **successfully learned from a video lecture on thin film interference** to correct its initial misunderstanding of phase shifts at boundaries.

The video demonstrates when phase shifts occur at interfaces between materials with different refractive indices. The adaptation question tests this knowledge by asking about the film thickness required for maximum constructive interference.

Before video, the model:

- Incorrectly assumed  $180^\circ$  phase shifts at both boundaries
- Failed to understand that phase shift depends on relative refractive indices
- Reached wrong conclusion that  $t = \lambda/2$  for constructive interference

After watching the video, the model:

- Correctly identified phase shift only at top boundary (low to high index)
- Applied proper interference conditions to determine  $t = \lambda/4$  for constructive interference

In conclusion, this case demonstrates successful knowledge acquisition from video. The model **corrects its understanding of phase shifts at interfaces** and thin film interference through learning from the video, and applies the demonstrated principles to solve the adaptation problem correctly.

Figure 20. A Wrong-to-Right example of Claude-3.5-Sonnet in the Adaptation track.



## Wrong-to-Right Example

### Video Lecture – "2-3 TREES - Search, Insertion and Deletion"

2-3 TREES

SHREEHAR JOSHI  
RAMAPO COLLEGE OF NEW JERSEY

Insertion (Case 1)

Insert in a node with only one data element.

Insert 4

Insertion (Case 2)

Insert in a node with two data elements whose parents contain only one data element.

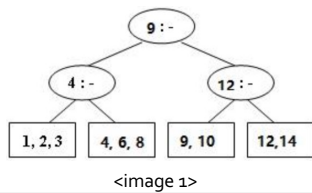
Insert 10

Insertion (Case 3)

Insert in a node with two data elements whose parents also contain two data elements.

Insert 1

### Adaptation Question



**Question:** After inserting 7 into the 2-3 tree given in the figure, how many of the following statements are TRUE?

- (S1) The root still contains 9 only;  
 (S2) the tree grows higher;  
 (S3) 4 and 9 are in the same nonleaf node;  
 (S4) the interior node containing 12 keeps unchanged. <image 1>

**Option:**

- (A) 3                      (B) 4                      (C) 0                      (D) 1                      **(E) 2** ✓

#### Before Video

##### Claude-3.5-Sonnet:

Let me help solve this step by step by analyzing what happens when we insert 7 into this 2-3 tree.  
 When inserting 7: 7 should go between 6 and 8 in the node containing 4,6,8

Let's analyze each statement:

- (S1) "The root still contains 9 only" - When we insert 7, the node (4,6,8) becomes too full. - This causes a split, pushing 6 up. - This affects the root, as now 6 and 9 must be in the root. - **Therefore FALSE.** ✗  
 (S2) "The tree grows higher" - The insertion causes internal reorganization. - The tree maintains its height. - **Therefore FALSE**  
 (S3) "4 and 9 are in the same nonleaf node" - After reorganization, 4 remains at its level - 9 moves to the right at its level - They are not in the same node - **Therefore FALSE.**  
 (S4) "The interior node containing 12 keeps unchanged"- The node containing 12 is not affected by inserting 7 - It remains unchanged - **Therefore TRUE.**

Counting TRUE statements: Only S4 is TRUE. All others (S1, S2, S3) are FALSE

Therefore, only 1 statement is TRUE. The answer is **D. 1** ✗

#### After Video

##### Claude-3.5-Sonnet:

Let me help solve this step by step by analyzing what happens when we insert 7 into the given 2-3 tree and checking each statement with what we learned from the video. When inserting 7:

- 7 would go between 6 and 8 in the leaf node [4,6,8]
- This would make that node overflow (too many elements)
- Following 2-3 tree insertion rules, **this would cause a split and reorganization.** This aligns with the principles of **Case 2 in the video** where the parent has 1 element. ✓

Let's check each statement:

- (S1) "The root still contains 9 only" - After insertion and reorganization, the root would still contain only 9 - **TRUE** ✓  
 (S2) "The tree grows higher" - The insertion of 7 doesn't cause the tree to grow in height. - The reorganization happens within the existing levels. - **FALSE**  
 (S3) "4 and 9 are in the same nonleaf node"- After reorganization, 4 and 9 would not be in the same node. - They remain separated at their current levels. - **FALSE**  
 (S4) "The interior node containing 12 keeps unchanged"- The right side of the tree (containing 12) is not affected by inserting 7 - That node remains unchanged- **TRUE**  
 Total number of TRUE statements = 2. Therefore, the answer is **E. 2.** ✓

#### Comment:

This example illustrates how the model successfully learned from a video lecture on 2-3 trees to **correct its initial misunderstanding of insertion operations and tree reorganization.**

The video demonstrates different cases of node insertion and the resulting reorganization rules in 2-3 trees. The adaptation question tests this knowledge by asking about the effects of inserting a new value into a given 2-3 tree.

Before video, the model:

- Incorrectly evaluated the effects of insertion on root node content
- Misunderstood how tree reorganization affects node relationships
- Wrongly concluded only 1 statement (S4) was true

After watching the video, the model:

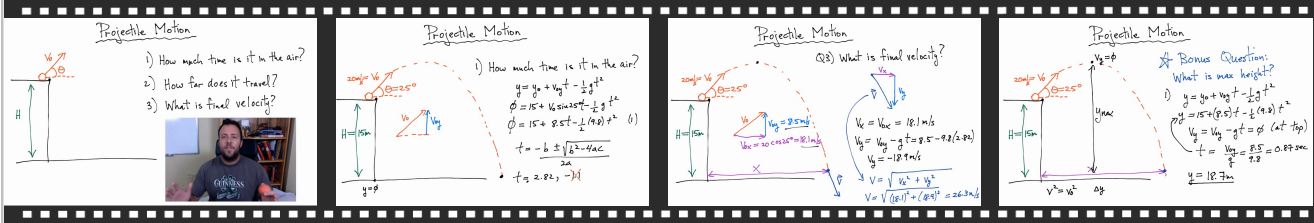
- Correctly identified how insertion causes node splits and reorganization
- Properly identify this scenario applies to the principles of Case 2 in the video
- Accurately determined 2 statements (S1 and S4) were true

In conclusion, this case demonstrates successful knowledge acquisition from video. The model corrects its understanding of 2-3 tree operations through learning from the video, and applies the demonstrated principles to solve the adaptation problem correctly.

Figure 21. A Wrong-to-Right example of Claude-3.5-Sonnet in the Adaptation track.

**Wrong-to-Right Example**

**Video Lecture – "Projectile launched off a cliff at an angle"**



**Adaptation Question**

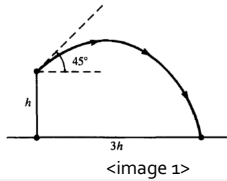


Fig. 19-1

**Question:** A rocket is shot from the top of a tower at an angle of  $45^\circ$  above the horizontal (Fig. 19-1). It hits the ground in 5 seconds at a horizontal distance from the foot of the tower equal to three times the height of the tower. Find the height of the tower. <image 1>

**Option:**

- (A)  **$h = 100$  ft** ✓
- (B)  $h = 80$  ft
- (C)  $h = 110$  ft
- (D)  $h = 85$  ft
- (E)  $h = 90$
- (F)  $h = 95$  ft
- (G)  $h = 105$  ft
- (H)  $h = 120$  ft
- (I)  $h = 75$  ft
- (J)  $h = 115$  ft

**Before Video**

**GPT-4o:**

To solve this problem, we need to use the principles of projectile motion.

Let's break it down step by step.

**1. Horizontal Motion:**

- The horizontal distance covered by the projectile is given by:

$$R = v_0 \times \cos(\theta) \times t$$

- Given  $R = 3h$ , we can write:

$$3h = v_0 \times \cos(45^\circ) \times 5$$

- Since  $\cos(45^\circ) = 1/\sqrt{2}$ , we have:

$$v_0 = (3h \times \sqrt{2})/5$$

**2. Vertical Motion:**

- The vertical displacement is given by:

$$h = v_0 \times \sin(45^\circ) \times t - (1/2) \times g \times t^2$$

- Since  $\sin(45^\circ) = 1/\sqrt{2}$ , we have:

$$h = [(3h \times \sqrt{2})/5 \times (1/\sqrt{2})] \times 5 - (1/2) \times 9.8 \times 5^2$$

$$h = [(3h)/5] \times 5 - (1/2) \times 9.8 \times 25$$

$$h = 3h - (1/2) \times 9.8 \times 25$$

$$h = 3h - 122.5$$

$$2h = 122.5$$

$$h = 61.25 \text{ ft}$$

However, this value does not match any of the given options. It seems there might be a mistake in the given options. The correct height based on the calculations is  $h = 61.25$  ft, which is not listed among the options. ✗

**After Video**

**GPT-4o:**

To solve the problem, we need to use the principles of projectile motion.

Here's a step-by-step solution:

**1. Horizontal Motion:**

- The horizontal distance (R) is given by:

$$R = v_0 \times \cos(\theta) \times t$$

- Given  $R = 3h$ :

$$3h = v_0 \times \cos(45^\circ) \times 5$$

- Since  $\cos(45^\circ) = 1/\sqrt{2}$ :

$$v_0 = (3h \times \sqrt{2})/5$$

**2. Vertical Motion:**

- The vertical displacement (y) is given by:

$$y = v_0 \times \sin(\theta) \times t - (1/2) \times g \times t^2$$

- **Since the rocket hits the ground,  $y = -h$ :** ✓

$$-h = v_0 \times \sin(45^\circ) \times 5 - (1/2) \times 9.8 \times 5^2$$

- Since  $\sin(45^\circ) = 1/\sqrt{2}$ :

$$-h = [(3h \times \sqrt{2})/5 \times (1/\sqrt{2})] \times 5 - (1/2) \times 9.8 \times 25$$

$$-h = 3h - (1/2) \times 9.8 \times 25$$

$$-h = 3h - 122.5$$

$$-4h = -122.5$$

$$h = 30.625 \text{ meters}$$

**3. Convert to Feet:**

- 1 meter = 3.28084 feet

- Therefore:

$$h = 30.625 \times 3.28084 = 100 \text{ feet. Answer: A. } h = 100 \text{ ft}$$
 ✓

**Comment:**

This example illustrates how the model successfully learned from a video lecture on projectile motion to **correct its initial errors**.

The video demonstrates an in-depth analysis of projectile motion launched from a cliff through four different question types: calculating time of flight, final velocity, horizontal distance, and maximum height. The adaptation question tests the application of this knowledge by asking to calculate the cliff height given the time of flight.

Before video, the model:

- Incorrectly used positive height ( $y = h$ ) in vertical motion equations, showing fundamental misunderstanding of direction
- Failed to convert between metric and imperial units properly

After watching the video, the model:

- Correctly used negative height ( $y = -h$ ) for downward displacement
- Properly handled unit conversions (meters to feet)
- And finally calculated the tower height as 100 feet

This case demonstrates successful knowledge acquisition from video. The model corrects its understanding of projectile motion principles through learning from the video and effectively applies the demonstrated concepts to solve the adaptation problem correctly.

Figure 22. A Wrong-to-Right example of GPT-4o in the Adaptation track.